

Université Paris I, Paris - Sorbonne

Première année Master M.A.E.F. 2007-2008

Statistiques

Plan du cours

1. Quelques rappels de la mesure.
2. Quelques rappels sur les applications de la mesure aux probabilités.
3. Estimation paramétrique.
4. Tests paramétriques.

Bibliographie

- Livres pour revoir les bases ...

1. Billard, B. Probabilités et techniques de l'analyse. SMG.
2. Berck, B., P. et Azo, E. Probabilités et applications - Cours Exercices. E. sciences.
3. Doss, F. Probabilités et Statistique. Dunod.
4. Lecut, J.-P. Statistiques et Probabilités.

Théorie de la mesure et applications aux probabilités

- Ansel et Duce, Exercices corrigés de la mesure et de l'intégration, Ellipses.
- Barbe, P. et Ledoux, M., Probabilités et Lin.
- Da Cunha - Caselle, D. et Duflo, M., Probabilités (I), Masson
- Jacod, J., Cours de Probabilités, <http://www.proba.jussieu.fr/pageperso/jacod.html>.
- Jacod, J., Cours de Probabilités, <http://www.proba.jussieu.fr/pageperso/jacod.html>.
- Toulouze, P. Théorie des probabilités et statistiques, Masson.

Statistiques inférentielles

- Da Cunha - Caselle, D. et Duflo, M., Probabilités (I), Masson.
- Fourdrinier, D., Statistiques inférentielles, Dunod.
- Lecut, J.-M. et Tassi, P., Statistiques et paramètres robustes, Economica.
- Milhaud, X., Statistiques, Belin.
- Monfort, A., Cours de statistiques inférentielles, Economica.
- Saporta, G., Probabilités et analyse des données et statistiques. Technip.
- Tsybakov, A. Introduction à la statistique et à l'analyse de données. Mathématiques et Applications, Springer.

Cours de STATISTIQUES

1 Rapports sur la théorie de la mesure

Introduction

Il demeure des choses inconnues à partir des connaissances probabilistes

- Quelles sont les propriétés et l'ensemble de tous les événements?
- Que se passe-t-il pour des probabilités d'événements moins classiques (par exemple l'ensemble des décimaux)?
- Comment tracer une variable aléatoire qui est continue et discrète à la fois (par exemple le nombre de minutes passées devant la TV)?

1.1 Mesures

1.1.1 Tribus

Notation. - Ω est un ensemble (fini ou infini).

- $\mathcal{P}(\Omega)$ est l'ensemble de tous les sous-ensembles (parties) de Ω .

Rappel. Soit E un ensemble. E est dit dénombrable s'il existe une bijection entre E ou un sous-ensemble de \mathbb{N} . Par exemple, un ensemble fini, $\mathbb{Z}, \mathbb{Z}^+, \mathbb{Z}^-$ sont dénombrables. En revanche, \mathbb{R} n'est pas dénombrable.

Définition. Soit une famille \mathcal{F} de parties de Ω (donc $\mathcal{F} \subset \mathcal{P}(\Omega)$). On dit que \mathcal{F} est une tribu si :

- $\Omega \in \mathcal{F}$;
- si $A \in \mathcal{F}$ alors $\Omega \setminus A \in \mathcal{F}$;
- pour tout $n \in \mathbb{N}^*$, si $(A_i)_{i \in \mathbb{N}^*} \in \mathcal{F}^{\mathbb{N}^*}$ alors $A_1 \cup \dots \cup A_n \in \mathcal{F}$.

Définition. Soit une famille \mathcal{A} de parties de Ω (donc $\mathcal{A} \subset \mathcal{P}(\Omega)$). On dit que \mathcal{A} est une σ -tribu (ou tribu σ) si :

- $\Omega \in \mathcal{A}$;
- si $A \in \mathcal{A}$ alors $\Omega \setminus A \in \mathcal{A}$;
- pour $I \subset \mathbb{N}$, si $(A_i)_{i \in I} \in \mathcal{A}^I$ alors $\bigcup_{i \in I} A_i \in \mathcal{A}$.

Exemple. - Cas du Pile ou Face.

- Cas où Ω est infini : $\Omega = \mathbb{N}$ par exemple.

Propriété. Avec les notations précédentes :

1. $\emptyset \in \mathcal{A}$;
2. si A et B sont dans la tribu, alors $A \cap B$ est dans \mathcal{A} ;
3. si A_1 et A_2 sont deux tribus sur Ω , alors $A_1 \cap A_2$ est une tribu sur Ω . Plus généralement, pour $I \subset \mathbb{N}$, si $(\mathcal{A}_i)_{i \in I}$ est une famille de tribus sur Ω , alors $\bigcap_{i \in I} \mathcal{A}_i$ est une tribu sur Ω ;
4. si A_1 et A_2 sont deux tribus sur Ω , alors $A_1 \cup A_2$ n'est pas forcément une tribu sur Ω . Par exemple, si $\Omega = \{0, 1, 2\}$, $\mathcal{T}_1 = \{\emptyset, \{0\}, \{1, 2\}, \Omega\}$ et $\mathcal{T}_2 = \{\emptyset, \{0, 1\}, \{2\}, \Omega\}$.

Définition. Si \mathcal{E} est une famille de parties de Ω (donc $\mathcal{E} \subset \mathcal{P}(\Omega)$), alors on appelle tribu engendrée par \mathcal{E} , notée $\sigma(\mathcal{E})$, la tribu engendrée par l'intersection de toutes les tribus \mathcal{F} contenant \mathcal{E} (on peut faire la même chose avec des σ -tribus).

Remarque. La tribu engendrée est la "plus petite" tribu (au sens de l'inclusion) contenant la famille \mathcal{E} .

Rappel. - Un ensemble ouvert U dans un espace métrique X est tel que pour tout $x \in U$, il existe $r > 0$ tel que $B(x, r) \subset U$.

- On dit qu'un ensemble dans un espace métrique X est fermé si son complémentaire dans X est ouvert.

Défini ti o n soit Ω un espace mesurable. On appelle tribu ou σ -algèbre sur Ω , noté $\mathcal{B}(\Omega)$, la tribu engendrée par les sous-ensembles de Ω . Un ensemble $B \in \mathcal{B}(\Omega)$ est appelé borélien.

Exemple. - Boréliens sur \mathbb{R} , sur $]0, 1[$. Attention $x - \frac{1}{10}; x + \frac{1}{10} \in B([0, 1])$ car une réunion d'ouverts est ouverte.
- Boréliens sur \mathbb{R}^2

1.1.2 Espace mesurable

Défini ti o n soit Ω un ensemble et \mathcal{A} une tribu sur Ω . On dit que (Ω, \mathcal{A}) est un espace mesurable.

Corollaire. Quand on s'intéresse aux probabilités, on dira que (Ω, \mathcal{A}) est un espace probabilisable.

Propriété. Si $(\Omega_i, \mathcal{A}_i)_i$ sont n espaces mesurables, alors un ensemble A de $\Omega = \Omega_1 \times \dots \times \Omega_n$ est une union finie d'ensembles $A_1 \times \dots \times A_n$ où chaque $A_i \in \mathcal{A}_i$. L'ensemble des ensembles mesurables est une algèbre et on note $\mathcal{A}_1 \otimes \dots \otimes \mathcal{A}_n$ (on dit \mathcal{A}_1 tensoriel $\mathcal{A}_2 \dots$ tensoriel \mathcal{A}_n) la tribu sur Ω engendrée par ces ensembles mesurables.

Exemple. Cas de \mathbb{R} .

Défini ti o n On appelle espace mesurable produit (Ω, \mathcal{A}) l'espace mesurable (Ω, \mathcal{A}) où $\mathcal{A} = \bigotimes_{i=1}^n \mathcal{A}_i$.

Exemple. Pile / Face 2 fois.

1.1.3 Définitions et Propriétés d'un espace mesurable

Défini ti o n soit (Ω, \mathcal{A}) un espace mesurable. L'application $\mu: [0, +\infty]$ est une mesure si:

- $\mu(\emptyset) = 0$.
- Pour tout $I \subset \mathbb{N}$ et pour $(A_i)_{i \in I}$ famille disjointe de \mathcal{A} (tel le que $A_i \cap A_j = \emptyset$ pour $i \neq j$), alors

$$\mu\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} \mu(A_i) \quad (\text{propriété dite de } \sigma\text{-additivité})$$

Défini ti o n Avec les notations précédentes:

- Si $\mu(\Omega) < +\infty$, on dit que μ est finie.
- Si $\mu(\Omega) < M$ avec $M < +\infty$, on dit que μ est bornée.
- Si $\mu(\Omega) = 1$, on dit que μ est une mesure de probabilité.

Exemple. Cas de $\Omega = \mathbb{R}$, de \mathbb{N} , ou \mathbb{R}^2

Défini ti o n Si (Ω, \mathcal{A}) est un espace mesurable (probabilisable) alors $(\Omega, \mathcal{A}, \mu)$ est un espace mesuré (resp. probabilisé) quand μ est une probabilité.

Remarque. Sur (Ω, \mathcal{A}) , on peut définir une infinité de mesures.

Propriété. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré et $(A_i)_{i \in \mathbb{N}}$ une famille de \mathcal{A} .

1. Si $A_1 \subset A_2$, alors $\mu(A_1) \leq \mu(A_2)$.
2. Si $\mu(A_1) < +\infty$ et $\mu(A_2) < +\infty$, alors $\mu(A_1 \cup A_2) + \mu(A_1 \cap A_2) = \mu(A_1) + \mu(A_2)$.

3. Pour tout $I \subset \mathbb{N}$, on a $\mu\left(\bigcup_{i \in I} A_i\right) \leq \sum_{i \in I} \mu(A_i)$.

4. Si $A_i \subset A_{i+1}$ pour tout $i \in \mathbb{N}$ (suite croissante en sens de l'inclusion), alors $(A_i)_{i \in \mathbb{N}}$ est une suite croissante convergente tel le que

$$\mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \lim_{i \rightarrow +\infty} \mu(A_i) \quad (\text{même si cette limite est } +\infty)$$

5. Si $A_{i+1} \subset A_i$ pour tout $i \in \mathbb{N}$ (suite décroissante en sens de l'inclusion) et $\mu(A_1) < +\infty$, alors

$$(\mu(A_i))_{i \in \mathbb{N}} \text{ est une suite décroissante convergente tel le que } \mu\left(\bigcap_{i \in \mathbb{N}} A_i\right) = \lim_{i \rightarrow +\infty} \mu(A_i).$$

Exemple 1. Soit (Ω, μ) un espace mesuré. On définit $v(A) = \mu(A)$ où $u \in A$. v mesure-t-elle ?

2. Si μ_1 et μ_2 mesures sur Ω , $\mu_1 + \mu_2$ et $\alpha\mu$ sont-elles des mesures ?

Définition 1. Soit (Ω, μ) un espace mesuré. Soit $(A_i)_{i \in \mathbb{N}}$ une famille de

1. On définit $\limsup_{n \rightarrow \infty} A_n = \bigcap_{n \in \mathbb{N}} \bigcup_{m \geq n} A_m$ (intuitivement, $\limsup_{n \rightarrow \infty} A_n$ est l'ensemble des $\omega \in \Omega$ tels que ω appartient à un nombre infini d'entre eux).

2. On définit $\liminf_{n \rightarrow \infty} A_n = \bigcup_{n \in \mathbb{N}} \bigcap_{m \geq n} A_m$ (intuitivement, $\liminf_{n \rightarrow \infty} A_n$ est l'ensemble des $\omega \in \Omega$ tels que ω appartient à tous les A_n à partir d'un certain rang).

Exemple 2. Cas des suites croissantes et décroissantes d'ensembles.

Théorème (Théorème d'extension de Hahn - Carathéodory). Soit ν une mesure sur Ω , et ν une application dans $[0, +\infty]$ additive (tel que $\nu(A \cup B) = \nu(A) + \nu(B)$ pour $A \cap B = \emptyset$), alors si A est la tribu engendrée par \mathcal{F} , il existe une mesure ν sur \mathcal{F} qui coïncide avec ν sur \mathcal{F} (c'est-à-dire que pour tout $E \in \mathcal{F}$, $\nu(E) = \nu(E)$). On dit que ν prolonge ν sur la tribu \mathcal{F} .

Exemple 3. Définition de la mesure de Lebesgue sur \mathbb{R} , \mathbb{R}^n , \mathbb{C} .

Définition 1. Soit (Ω, μ) un espace mesuré.

- Pour $A \in \mathcal{A}$, on dit que A est μ -négligeable si $\mu(A) = 0$.
- Soit une propriété P dépendant des éléments ω de Ω . On dit que P est vraie μ -presque partout (ou μ -presque sûrement) sur un espace mesuré si l'ensemble des ω pour lequel elle n'est pas vraie est μ -négligeable.

Exemple 4. - Mesure de Lebesgue sur \mathbb{R} , \mathbb{R}^n , \mathbb{C} .

- La propriété "la suite de fonctions $f_n(x) = x^n$ converge vers la fonction $f(x) = 0$ " est vraie λ -presque partout sur $[0, 1]$.
- Soit $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mu)$ et soit F la fonction définie par $F(x) = \mu(\{y \in \mathbb{R}, y \leq x\})$ pour $x \in \mathbb{R}$.

1.1.4 Fonctions mesurables

Rappel. Soit $f: E \rightarrow F$, où E et F sont 2 espaces métriques.

- Pour $I \subset F$, on appelle l'ensemble préimage de I par f , l'ensemble $f^{-1}(I) = \{x \in E, f(x) \in I\}$.
- $(f$ continue) \Leftrightarrow (pour tout ouvert U de F alors $f^{-1}(U)$ est un ouvert de E).

Définition 1. Soit $f: E \rightarrow F$ et soit I une tribu sur F . On note $\sigma(f)$ l'ensemble des sous-ensembles de E tel que $f^{-1}(I) \in \mathcal{A}$ (où $I \in \mathcal{I}$).

Propriété. Soit (Ω, \mathcal{A}) un espace mesuré et soit $f: \Omega \rightarrow F$. Alors $\sigma(f^{-1}(\mathcal{I}))$ est une tribu sur Ω engendrée par f .

Définition 2. Soit (Ω, \mathcal{A}) et (Ω', \mathcal{A}') deux espaces mesurables. Une fonction $f: \Omega \rightarrow \Omega'$ est dite mesurable pour les tribus \mathcal{A} et \mathcal{A}' si et seulement si $f^{-1}(A') \in \mathcal{A}$ (donc si et seulement si $f^{-1}(A') \in \mathcal{A}$, alors $\sigma(f^{-1}(\mathcal{A}')) \in \mathcal{A}$).

Exemple 5. - Fonction indicatrice.

- Combinons les tribus de fonctions indicatrices.

Remarque 1. Dans le cas où $\Omega = \mathbb{R}$, $\mathcal{A} = \mathcal{B}(\mathbb{R})$ est un espace probabilisable, et si $f: \mathbb{R} \rightarrow \mathbb{R}$ alors si f est une fonction mesurable sur $\mathcal{B}(\mathbb{R})$, alors f est une variable aléatoire.

Exemple 6. Nombre de Piles dans un jeu de Pile/Face.

Remarque 2. Dans le cas où $\Omega = \{0, 1\}^n$ est un espace mesurable, et si $f: \Omega \rightarrow \mathbb{R}$, où Ω est un espace métrique et $\mathcal{B}(\Omega)$ l'ensemble des boréliens de Ω si f est une fonction mesurable sur $\mathcal{B}(\Omega)$, alors f est dite fonction borélienne.

Proposition 1. Soit (Ω, \mathcal{A}) et (Ω', \mathcal{A}') deux espaces mesurables et $f: \Omega \rightarrow \Omega'$. Soit \mathcal{F} une famille de sous-ensembles de Ω telle que $\bigcup \mathcal{F} = \Omega$. Alors

1. $f^{-1}(F)$ engendre la tribu $\mathcal{F}(A)$.

2. $(f \text{ mesurable}) \Rightarrow (f^{-1}(F) \subset A)$

Conséquence. - Si (Ω, A) et (Ω', A') sont deux espaces mesurables, alors toute application continue de Ω dans Ω' est mesurable.

- Pour montrer qu'une fonction $f: \Omega \rightarrow \mathbb{R}$ est mesurable, il suffit de montrer que la famille d'ensemble $(\{\omega \in \Omega, f(\omega) \leq a\})_{a \in \mathbb{R}} \in A$.

Propriété. - Soit f mesurable de (Ω, A) dans (Ω', A') et g mesurable de (Ω, A) dans (Ω'', A'') . Alors $g \circ f$ est mesurable de (Ω, A) dans (Ω'', A'') .

- Soit f_1 mesurable de (Ω, A) dans (Ω_1, A_1) et f_2 mesurable de (Ω, A) dans (Ω_2, A_2) . Alors $h: \Omega \rightarrow \Omega_1 \times \Omega_2$ tel que $h(\omega) = (f_1(\omega), f_2(\omega))$ est mesurable de (Ω, A) dans $A_1 \otimes A_2$.

- Soit $(f_n)_{n \in \mathbb{N}}$ une suite de fonctions mesurables (Ω, A) dans $(\Omega', B(\Omega'))$, où Ω' est un espace métrique tel qu'il existe une fonction f limite simple de (f_n) (c'est-à-dire $\forall \omega \in \Omega, \lim_{n \rightarrow \infty} f_n(\omega) = f(\omega)$).

Alors f est mesurable de (Ω, A) dans $(\Omega', B(\Omega'))$.

Définition. Soit f mesurable de (Ω, A) dans (Ω', A') et soit $\mu: A' \rightarrow [0, +\infty]$ tel que pour tout $A' \in A'$, on ait $\mu_f(A') = \mu(f^{-1}(A'))$. Alors μ_f est une mesure sur A appelée mesure image de μ par f .

Cas particuliers. Si μ est une mesure de probabilité et si f est une variable aléatoire, alors μ_f est la mesure (loi) de probabilité de la variable aléatoire X .

1.1.5 Cas des fonctions réelles mesurables

Propriété. Soit f et g deux fonctions réelles mesurables (de (Ω, A) dans $(\mathbb{R}, B(\mathbb{R}))$). Alors αf , $f + g$, $\min(f, g)$ et $\max(f, g)$ sont des fonctions réelles mesurables.

Propriété. Soit $(f_n)_{n \in \mathbb{N}}$ une suite de fonctions réelles mesurables. Alors $\inf_n f_n$ et $\sup_n f_n$ sont des fonctions réelles mesurables.

Définition. Soit $f: \Omega \rightarrow \mathbb{R}$. Alors f est dite *étalée* s'il existe une famille d'ensembles disjoints $(A_i)_{i \in \mathbb{N}}$ de Ω et une famille de réels $(\alpha_i)_{1 \leq i \leq n}$ tel que pour tout $\omega \in \Omega$, on ait $f(\omega) = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}(\omega)$.

Remarque. Si les A_i sont tous dans A et si Ω est distribué sur Ω , alors f est mesurable.

Théorème. Toute fonction réelle mesurable à valeurs dans \mathbb{R} est limite simple d'une suite croissante de fonctions étalées.

Conséquence. Soit f une fonction réelle mesurable. Alors f est limite simple de fonctions étalées.

1.2 Intégration de Lebesgue

Dans toute la suite, on considère (Ω, A, μ) un espace mesuré.

1.2.1 Intégrale de Lebesgue d'une fonction positive

Définition 1. Soit $f = \mathbb{1}_A$, où $A \in A$. Alors :

$$\int_{\Omega} f d\mu = \int_{\Omega} \mathbb{1}_A(\omega) d\mu(\omega) = \mu(A).$$

2. Soit $f = \mathbb{1}_B$, où $B \in A$ et soit $A \in A$. Alors :

$$\int_B f d\mu = \int_B \mathbb{1}_A(\omega) d\mu(\omega) = \mu(A \cap B) = \mu(B).$$

3. Soit f une fonction étalée positive tel que $f = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}$, où les $A_i \in A$ et $\alpha_i > 0$ et soit $B \in A$.

Alors :

$$\int_B f d\mu = \int_B \left(\sum_{i=1}^n \alpha_i \mathbb{1}_{A_i}(\omega) \right) d\mu(\omega) = \sum_{i=1}^n \alpha_i \mu(A_i \cap B).$$

Exemple : Fonction χ_A , fonction en escalier,...

Définition Soit f une fonction μ -mesurable positive et soit $A \subset B$. Alors l'intégrale de Lebesgue de f par rapport à μ sur B est :

$$\int_B f d\mu = \int_B \chi_B(\omega) f(\omega) d\mu(\omega) = \sup_B \int_B g d\mu, \text{ pour } g \text{ fonction positive telle que } g \leq f.$$

Propriété. Soit f une fonction μ -mesurable positive et soit $A \subset B$. Alors :

1. Pour $c \geq 0$, $\int_B cf d\mu = c \int_B f d\mu$.
2. Si $A \subset B$, alors $\int_A f d\mu \leq \int_B f d\mu$.
3. Si g est une fonction μ -mesurable positive telle que $0 \leq g \leq f$ alors $0 \leq \int_B g d\mu \leq \int_B f d\mu$.
4. Si $\mu(B) = 0$ alors $\int_B f d\mu = 0$.

Théorème (Théorème de convergence monotone (Beppo Levi)) Soit $(f_n)_n$ est une suite croissante de fonctions mesurables positives convergeant simplement vers f sur Ω , alors :

$$\lim_{n \rightarrow \infty} \int_{\Omega} f_n d\mu = \int_{\Omega} f d\mu = \lim_{n \rightarrow \infty} \int_{\Omega} f_n d\mu.$$

Cependant pour les séries de fonctions mesurables positives, on peut toujours appliquer le théorème de convergence monotone et donc inverser la somme et l'intégrale.

Lemme (Lemme de Fatou) Soit $(f_n)_n$ est une suite de fonctions mesurables positives alors :

$$\liminf_{n \rightarrow \infty} \int_{\Omega} f_n d\mu \leq \int_{\Omega} \liminf_{n \rightarrow \infty} f_n d\mu.$$

Exemple Appliquer Fatou à $f_n = \chi_A$ et $f_{2n+1} = \chi_B$.

1.2.2 Intégrale de Lebesgue d'une fonction réelle et propriétés

Définition Soit (Ω, μ) un espace mesuré et soit f une fonction μ -mesurable à valeurs réelles tel que $f = f^+ - f^-$ avec $f^+ = \max(f, 0)$ et $f^- = \max(-f, 0)$. On dit que f est μ -intégrable sur B si

$$\int_B |f| d\mu < +\infty. \text{ On a alors}$$

$$\int_B f d\mu = \int_B f^+ d\mu - \int_B f^- d\mu.$$

Notation Lorsque f est μ -intégrable sur B , soit $\int_B |f| d\mu < +\infty$, on note $f \in L^1(\Omega, \mu)$ (on dit que f est L^1).

Exemple L'intégrale de Riemann est une intégrale de Lebesgue.
Cas de la mesure de Dirac.

Propriété. On suppose que $f \in L^1(\Omega, \mu)$. Alors :

1. $(\alpha f + \beta g) d\mu = \alpha \int f d\mu + \beta \int g d\mu$ pour $(\alpha, \beta) \in \mathbb{R}^2$.
2. Si $f \leq g$ alors $\int f d\mu \leq \int g d\mu$.

Théorème (Théorème de convergence dominée de Lebesgue) Soit $(f_n)_n$ est une suite de fonctions de $L^1(\Omega, \mu)$ telles que pour tout $n \in \mathbb{N}$, $|f_n| \leq g$ avec $g \in L^1(\Omega, \mu)$. Si on suppose que f_n converge simplement vers f sur Ω alors :

$$\lim_{n \rightarrow \infty} \int_{\Omega} f_n d\mu = \int_{\Omega} f d\mu.$$

Exemple 1. Soit $f: \mathbb{R} \rightarrow \mathbb{R}$ une fonction convexe et soit $f: \Omega \rightarrow \mathbb{R}$ mesurable tel que $\varphi(f)$ soit une fonction intégrable par rapport à P . Alors:

$$\varphi \int f dP \leq \int \varphi(f) dP.$$

Exemple 2. Soit X une v.a. sur (Ω, \mathcal{F}, P) . Alors $\varphi(\int X dP) \leq \int \varphi(X) dP$.

1.2.3 Mesures induites et densité

Théorème 1. Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré et f une fonction mesurable positive. Soit ν la mesure induite par f (donc $\nu(A) = \int_A f d\mu$) pour $A \in \mathcal{F}$ et soit φ une fonction mesurable de (Ω, \mathcal{F}) dans $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$. Alors, si $\varphi \in L^1(\Omega, \nu)$,

$$\int_{\Omega} \varphi d\nu = \int_{\Omega} \varphi f d\mu.$$

Définition 1. Soit μ et ν deux mesures sur (Ω, \mathcal{F}) . On dit que μ domine ν (ou ν est dominé par μ) et que ν est absolument continu par rapport à μ lorsque $\mu(A) = 0 \Rightarrow \nu(A) = 0$.

Propriété. Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré et f une fonction positive mesurable et positive. On suppose que pour $A \in \mathcal{F}$, $\nu(A) = \int_A f d\mu$. Alors, ν est une mesure sur (Ω, \mathcal{F}) dominée par μ . De plus, pour toute fonction positive mesurable et positive,

$$\int g d\nu = \int g f d\mu.$$

Enfin, g est ν -intégrable si et seulement si $g f$ est μ -intégrable.

Définition 2. On dit que μ mesure sur (Ω, \mathcal{F}) est σ -finie lorsque existe une famille $\{A_i\}_{i \in I}$, avec I dénombrable, d'ensembles tels que $A_i = \Omega$ et $\mu(A_i) < +\infty$ pour tout $i \in I$.

Théorème 2 (Théorème de Radon-Nikodym). Soit μ et ν deux mesures σ -finies sur (Ω, \mathcal{F}) , telles que μ domine ν . Alors, existe une fonction positive mesurable et positive, appelée densité de ν par rapport à μ , tel que pour tout $A \in \mathcal{F}$, $\nu(A) = \int_A f d\mu$.

Théorème 3 (Théorème de Fubini). Soit $\Omega = \Omega_1 \times \Omega_2$, $\mathcal{F} = \mathcal{F}_1 \otimes \mathcal{F}_2$ et $\mu = \mu_1 \otimes \mu_2$ (mesures σ -finies), où $(\Omega_1, \mathcal{F}_1, \mu_1)$ et $(\Omega_2, \mathcal{F}_2, \mu_2)$ sont des espaces mesurés. Soit une fonction $f: \Omega \rightarrow \mathbb{R}$, \mathcal{F} -mesurable et μ -intégrable. alors:

$$\int_{\Omega} f d\mu = \int_{\Omega_1} \left(\int_{\Omega_2} f(\omega_1, \omega_2) d\mu_2(\omega_2) \right) d\mu_1(\omega_1) = \int_{\Omega_2} \left(\int_{\Omega_1} f(\omega_1, \omega_2) d\mu_1(\omega_1) \right) d\mu_2(\omega_2).$$

1.2.4 Espaces L^p

Définition 3. Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré. On appelle espace $L^p(\Omega, \mathcal{F}, \mu)$, où $p > 0$, l'ensemble des fonctions $f: \Omega \rightarrow \mathbb{R}$, mesurables et telles que $\int |f|^p d\mu < +\infty$.

Définition 4. Pour $f \in L^p(\Omega, \mathcal{F}, \mu)$, où $p > 0$, on note $\|f\|_p = \left(\int |f|^p d\mu \right)^{1/p}$.

Propriété (Inégalité Hölder). Soit $p > 1$ et $q > 1$ tels que $\frac{1}{p} + \frac{1}{q} = 1$, et soit $f \in L^p(\Omega, \mathcal{F}, \mu)$ et $g \in L^q(\Omega, \mathcal{F}, \mu)$. Alors, $fg \in L^1(\Omega, \mathcal{F}, \mu)$ et

$$\|fg\|_1 \leq \|f\|_p \cdot \|g\|_q.$$

Propriété (Inégalité de Minkowski) Soit $p > 1$ et soit $f, g \in L^p(\Omega, A, \mu)$. Alors, $f + g \in L^p(\Omega, A, \mu)$ et

$$\|f + g\|_p \leq \|f\|_p + \|g\|_p.$$

Remarque. Pour $p > 1$, L^p définit un espace normé. Pour obtenir une norme, il faut se placer dans l'espace $L^p(\Omega, A, \mu)$ obtenu en "quotientant" $L^p(\Omega, A, \mu)$ par la relation d'équivalence $f \sim g$ μ -presque partout (c'est-à-dire que dans $L^p(\Omega, A, \mu)$ on dira que $f = g$ lorsque $f = g$ μ -presque partout).

Définition. Pour $f, g \in L^2(\Omega, A, \mu)$, on définit le produit scalaire $\langle f, g \rangle = \int f \bar{g} d\mu$. On munira ainsi $L^2(\Omega, A, \mu)$ d'une structure d'espace de Hilbert. On dira que f est orthogonal à g lorsque $\langle f, g \rangle = 0$.

Conséquence. Si A est un sous-espace vectoriel fermé de $L^2(\Omega, A, \mu)$ (par exemple un sous-espace de dimension finie), alors pour tout $f \in L^2(\Omega, A, \mu)$, il existe un unique projeté orthogonal de f sur A , noté f_A , qui vérifie $f_A = \text{Arg inf}_{g \in A} \|f - g\|_2$.

2 Applications de la théorie de la mesure et de l'intégration en Probabilités

2.1 Espérance de variables aléatoires

Définition. Soit X une variable aléatoire sur $(\Omega, \mathcal{F}, \mathbb{P})$ un espace probabilisé. Alors si $X \in L^1(\Omega, A, \mathbb{P})$, on définit l'espérance de X par le nombre $\mathbb{E}X = \int X d\mathbb{P}$. Plus généralement, si $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ est bornée et si $\varphi(X) \in L^1(\Omega, A, \mathbb{P})$, on définit l'espérance de $\varphi(X)$ par $\mathbb{E}\varphi(X) = \int \varphi(X) d\mathbb{P}$.

Propriété. Si X est une variable aléatoire sur $(\Omega, \mathcal{F}, \mathbb{P})$, si $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ est bornée et si $\varphi(X) \in L^1(\Omega, A, \mathbb{P})$, et si \mathbb{P}_X est la mesure de probabilité de X alors :

$$\mathbb{E}\varphi(X) = \int_{\mathbb{R}} \varphi(x) d\mathbb{P}_X(x).$$

Conséquence. - Si \mathbb{P}_X est absolument continue par rapport à la mesure de Lebesgue (donc X est une v.a. dite absolument continue), on a $\mathbb{E}\varphi(X) = \int_{\mathbb{R}} \varphi(x) f_X(x) dx$.

- Si \mathbb{P}_X est absolument continue par rapport à la mesure de comptage sur \mathbb{N} (donc X est une v.a. dite discrète), de densité p_X , alors $\mathbb{E}\varphi(X) = \sum_{k=0}^{\infty} p_X(k) \varphi(k)$.

Propriété. 1. Soit X et Y des variables aléatoires telles que $X, Y \in L^1(\Omega, A, \mathbb{P})$. Alors pour tout $(a, b) \in \mathbb{R}^2$, $aX + bY \in L^1(\Omega, A, \mathbb{P})$ et

$$\mathbb{E}(aX + bY) = a\mathbb{E}X + b\mathbb{E}Y.$$

2. Soit X une variable aléatoire sur $(\Omega, \mathcal{F}, \mathbb{P})$, et soit $A \in \mathcal{F}$. Alors $\mathbb{E}(\mathbb{1}_A(X)) = \mathbb{P}(X \in A)$.

3. Soit X et Y des variables aléatoires telles que $X \in L^p(\Omega, A, \mathbb{P})$ et $Y \in L^q(\Omega, A, \mathbb{P})$ avec $\frac{1}{p} + \frac{1}{q} = 1$ et $p > 1, q > 1$. Alors $XY \in L^1(\Omega, A, \mathbb{P})$ et

$$\mathbb{E}|X \cdot Y| \leq (\mathbb{E}|X|^p)^{1/p} (\mathbb{E}|Y|^q)^{1/q}.$$

4. Soit X et Y des variables aléatoires telles que $X, Y \in L^p(\Omega, A, \mathbb{P})$, avec $p \geq 1$. Alors $X + Y \in L^p(\Omega, A, \mathbb{P})$ et

$$(\mathbb{E}|X + Y|^p)^{1/p} \leq (\mathbb{E}|X|^p)^{1/p} + (\mathbb{E}|Y|^p)^{1/p}.$$

5. Soit X une variable aléatoire telle que $X \in L^p(\Omega, A, \mathbb{P})$ pour $p > 0$. Alors pour tout $0 < r \leq p$, $X \in L^r(\Omega, A, \mathbb{P})$ et

$$(\mathbb{E}|X|^r)^{1/r} \leq (\mathbb{E}|X|^p)^{1/p}.$$

6. Si X est une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, si $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ est une fonction réelle convexe telle que X et $\varphi(X) \in L^1(\Omega, \mathcal{A}, \mathbb{P})$, alors

$$\mathbb{E}(\varphi(X)) \geq \varphi(\mathbb{E}X).$$

Définition. Pour X et Y des variables aléatoires telles que $X, Y \in L^2(\Omega, \mathcal{A}, \mathbb{P})$, on définit la covariance de X et Y par

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)];$$

On appelle variance de X , $\text{var}(X) = \text{cov}(X, X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}(X^2) - (\mathbb{E}X)^2$.

Propriété. Sur $L^2(\Omega, \mathcal{A}, \mathbb{P})$, $\text{cov}(\cdot, \cdot)$ définit un produit scalaire. De plus

$$|\text{cov}(X, Y)|^2 \leq \text{var}(X) \cdot \text{var}(Y).$$

2.2 Fonction de répartition et quantiles d'une loi de probabilité

Il y a une correspondance bijective entre la fonction de répartition $F_X :]-\infty, x] \rightarrow [0, 1]$ et la fonction de densité $f_X : \mathbb{R} \rightarrow [0, +\infty]$. La fonction de répartition permet de définir les quantiles qui sont essentiels à la construction d'intervalle de confiance et de test.

Soit $\alpha \in [0, 1]$. Des propriétés de la fonction de répartition, on déduit qu'il existe $x_\alpha \in \mathbb{R}$, tel que :

$$\lim_{x \rightarrow x_\alpha^-} F_X(x) \leq \alpha \leq F_X(x_\alpha). \quad (1)$$

Soit $I_\alpha = \{x_\alpha \in \mathbb{R} \text{ tel que } x_\alpha \text{ vérifie (1)}\}$. On appelle quantile (ou fractile, ou percentile en anglais) d'ordre α de la loi \mathbb{P} , noté q_α , le milieu de l'intervalle I_α . Evidemment, lorsque X admet une distribution absolument continue par rapport à la mesure de Lebesgue, $q_\alpha = F_X^{-1}(\alpha)$, où F_X^{-1} désigne la fonction inverse de F_X .

Deux cas particuliers sont très importants :

1 / pour $\alpha = 0.5$, c'est appelé médiane de \mathbb{P} ;

2 / pour $\alpha = 0.25$ et $\alpha = 0.75$ (respectivement), q_α sont appelés premier et troisième quantile (respectivement) de \mathbb{P} .

3 / pour $\alpha = 0.1, \dots, 0.9$, on parle de déciles de \mathbb{P} .

2.3 Principales lois de probabilité

Loi uniforme discrète :

C'est la loi de probabilité discrète à valeurs dans $\{x_1, \dots, x_n\}$ telle que

$$\mathbb{P}(X = x_i) = \frac{1}{n}.$$

On a alors : $\mathbb{E}X = \frac{1}{n}(x_1 + \dots + x_n)$ et $\text{var}(X) = \frac{1}{n}(x_1^2 + \dots + x_n^2) - (\mathbb{E}X)^2$.

Loi de Bernoulli :

C'est la loi de probabilité discrète notée $B(p)$ à valeurs dans $\{0, 1\}$ telle que

$$\mathbb{P}(X = 1) = p \quad \text{et} \quad \mathbb{P}(X = 0) = 1 - p.$$

On a alors : $\mathbb{E}X = p$ et $\text{var}(X) = p(1-p)$.

Loi binomiale :

C'est la loi de probabilité discrète notée $B(n, p)$ à valeurs dans $\{0, 1, \dots, n\}$ telle que

$$\mathbb{P}(X = k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k} \quad \text{pour } k \in \{0, 1, \dots, n\}$$

On a alors : $X = X_1 + \dots + X_n$, où (X_i) est une suite de v.a. i.i.d. ~~Bé~~ $\text{Exp}(\theta)$ et $\text{var}(X) = n\theta(1-p)$.

Loi de Poisson :

C'est la loi de probabilité à valeurs dans \mathbb{N} telle que

$$\mathbb{P}(X = k) = \frac{\theta^k}{k!} \cdot e^{-\theta} \quad \text{pour } k \in \mathbb{N}.$$

On a alors $\mathbb{E}X = \theta$ et $\text{var}(X) = \theta$.

Loi uniforme sur $[a, b]$

Cette loi est générée par $U([a, b])$, $-\infty < a < b < \infty$. C'est la loi de probabilité à valeurs dans $[a, b]$ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{b-a} \cdot \mathbb{1}_{x \in [a, b]}$$

On a alors $\mathbb{E}X = \frac{b+a}{2}$ et $\text{var}(X) = \frac{(b-a)^2}{12}$.

Loi Gamma :

Cette loi est générée par $\gamma(p, \theta)$, $p > 0$ et $\theta > 0$. C'est la loi de probabilité à valeurs dans \mathbb{R}_+ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{\theta^p}{\Gamma(p)} \cdot e^{-\theta \cdot x} \cdot x^{p-1} \cdot \mathbb{1}_{x \in \mathbb{R}_+}.$$

On a alors $\mathbb{E}X = \frac{p}{\theta}$ et $\text{var}(X) = \frac{p}{\theta^2}$.

Si $X \sim \gamma(p, \theta)$ et $Y \sim \gamma(q, \theta)$ avec X et Y indépendantes et $p > 0$ et $q > 0$, alors $X+Y \sim \gamma(p+q, \theta)$.
Pour $p = 1$, la loi $\gamma(p, \theta)$ est la loi exponentielle.

Loi Béta :

Cette loi est générée par $\beta(p, q)$, $p > 0$ et $q > 0$. C'est la loi de probabilité à valeurs dans $[0, 1]$ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{x^p(1-x)^{q-1}}{B(p, q)} \cdot \mathbb{1}_{x \in [0, 1]} \quad \text{où } B(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$$

On a alors $\mathbb{E}X = \frac{B(p+1, q)}{B(p, q)}$ et $\text{var}(X) = \frac{p \cdot q}{(p+q)(p+q+1)}$

Si $X \sim \gamma(p, \theta)$ et $Y \sim \gamma(q, \theta)$ avec X et Y indépendantes et $p > 0$ et $q > 0$, alors $\frac{X}{X+Y} \sim \beta(p, q)$.
Pour $p = 1$, la loi $\gamma(p, \theta)$ est la loi exponentielle.

Loi normale (ou gaussienne) centrée :

Cette loi est générée par $\mathcal{N}(0, 1)$. C'est la loi de probabilité à valeurs dans \mathbb{R} de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right).$$

On a :

$$\mathbb{E}(X) = 0 \quad \text{et} \quad \text{var}(X) = 1.$$

Loi normale (ou gaussienne) de moyenne m et de variance σ^2

Si Z suit la loi $N(0, 1)$, $X = m + \sigma Z$ suit par définition la loi $N(m, \sigma^2)$, loi normale d'espérance m et de variance σ^2 . La densité de X est donc donnée par :

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x-m)^2}{2\sigma^2}.$$

La figure A.1. représente la densité de la loi normale centrée réduite et celle d'une loi normale non centrée et non réduite. A partir de la loi gaussienne, on peut en déduire les lois suivantes.

Loi du χ^2 à n degrés de liberté :

Soit X_1, \dots, X_n , n variables aléatoires indépendantes de $N(0, 1)$, alors

$$S = X_1^2 + \dots + X_n^2$$

suit une loi du χ^2 à n degrés de liberté, loi notée $\chi^2(n)$. Cette loi est à valeurs dans \mathbb{R}_+ et sa densité est donnée par la loi Gamma $\gamma(n/2, 1/2)$, c'est-à-dire $\exp(-x/2) x^{n/2-1}$ pour $x \geq 0$ et 0 ailleurs. On peut aussi la caractériser par son rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{2^{n/2} \cdot \Gamma(n/2)} x^{n/2-1} \exp -\frac{x}{2} \cdot \mathbb{I}_{\{x \geq 0\}}$$

où la fonction Gamma est telle que $\Gamma(a) = \int_0^\infty x^{a-1} \cdot e^{-x} dx$ pour $a \geq 0$. Enfin, si X suit une loi $\chi^2(n)$, par définition on dira que $Y = \sqrt{X}$ suit une loi $\chi(n)$. La figure A.2. exhibe trois exemples de densités de loi du χ^2 . Loi de Student à n degrés de liberté :

La loi de Student à n degrés de liberté notée $T(n)$, est la loi du quotient

$$T = \frac{N}{S/n}$$

où N suit une loi $N(0, 1)$ et S suit une loi $\chi^2(n)$, N et S étant deux variables aléatoires indépendantes. On peut également caractériser la densité d'une telle loi par rapport à la mesure de Lebesgue, à savoir,

$$f_X(x) = \frac{1}{\sqrt{n} \cdot B(1/2, n/2)} \left(1 + \frac{x^2}{n}\right)^{-\frac{(n+1)}{2}},$$

où la fonction Beta est telle que $B(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)}$ pour $a > 0$ et $b > 0$. La figure A.3 illustre deux exemples de cette densité l'on compare à l'élément avec la densité loi normale centrée réduite.

Remarque : Pour la loi des grands nombres, plus n est grand, plus S est proche de sa valeur espérée. Le dénominateur est donc proche de 1. Il s'ensuit que la loi $T(n)$ est d'autant plus proche d'une loi normale que n est grand.

Un des principaux intérêts de la loi de Student est de nous permettre de caractériser la loi de X si X_1, X_n sont n variables aléatoires indépendantes de $N(m, \sigma^2)$, si on connaît la moyenne et la variance empiriques :

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n) \text{ et } \bar{S}_n^2 = \frac{1}{n-1} (X_1 - \bar{X}_n)^2 + \dots + (X_n - \bar{X}_n)^2,$$

alors

$$T = \frac{\sqrt{n} \cdot (\bar{X}_n - m)}{\bar{S}_n}$$

suit une loi de Student à n degrés de liberté.

Loi de Fisher à 1 et 2 degrés de liberté :

Soit S_1 et S_2 deux variables aléatoires indépendantes de lois respectives $\chi^2(n_1)$ et $\chi^2(n_2)$. Alors par définition :

$$F = \frac{S_1/n_1}{S_2/n_2}$$

suit une loi de Fisher à n_1 et n_2 degrés de liberté notée $F(n_1, n_2)$.

Remarque : Par les mêmes considérations que précédemment, la loi F est d'autant plus proche de 1 que les degrés de liberté n_1 et n_2 sont grands.

On a également les propriétés suivantes :

- Si F suit une loi $F(n_1, n_2)$, alors la loi de $\frac{n_1}{n_2} F$ est une loi beta de seconde espèce de paramètres $(n_1/2, n_2/2)$, c'est-à-dire que F est à valeur standard et la densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{B(n_1/2, n_2/2)} n_1^{n_1/2} \cdot n_2^{n_2/2} \frac{x^{n_1/2-1}}{(n_2 + n_1 \cdot x)^{(n_1+n_2)/2}} \mathbb{I}_{\{x \geq 0\}}$$

la notation Beta n'ayant encore la fonction Beta.

- Si $F \sim F(n_1, n_2)$, alors $E(F) = \frac{n_2}{n_2 - 2}$ lorsque $n_2 > 2$ et $\text{var}(F) = \frac{2n_2^2(n_1 + n_2 - 2)}{n_1(n_2 - 4)(n_2 - 2)^2}$ lorsque $n_2 > 4$.
- Si T suit une loi de Student $T(n)$, alors T^2 suit une loi de Fisher $F(1, n)$.

La figure A.4. donne une idée de la distribution d'une loi de Fisher pour différents choix des paramètres.

2.4 Indépendance

Définitions : Soit $(\Omega, \mathcal{F}, \mathbb{P})$ un espace probabilisé :

- Soit $(A_i)_{i \in I}$ une famille d'événements. On dit que les événements $(A_i)_{i \in I}$ sont indépendants si et seulement si pour tous les sous-ensembles finis K

$$\mathbb{P} \left(\bigcap_{i \in K} A_i \right) = \prod_{i \in K} \mathbb{P}(A_i).$$

- Soit $(A_i)_{i \in I}$ une famille de sous-tribus. On dit donc pour tout $i \in I$, $A_i \subset \mathcal{A}$. On dit que les tribus $(A_i)_{i \in I}$ sont indépendantes si et seulement si pour tous les sous-ensembles finis K les événements $\bigcap_{i \in K} A_i$ avec $K \subset I$ sont indépendants.
- Soit $(X_i)_{i \in I}$ des variables aléatoires sur $(\Omega, \mathcal{F}, \mathbb{P})$, à valeurs dans $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. On dit que les v.a. $(X_i)_{i \in I}$ sont indépendantes si et seulement si les tribus $\sigma(X_i)_{i \in I}$ sont indépendantes.

Proposition : Si (X_1, \dots, X_n) sont des variables aléatoires sur $(\Omega, \mathcal{F}, \mathbb{P})$. Alors les (X_i) sont indépendantes

si et seulement si $\mathbb{P}_{(X_1, \dots, X_n)} = \prod_{i=1}^n \mathbb{P}_{X_i}$.

Proposition : Si $(X_i)_{i \in I}$ sont des variables aléatoires indépendantes sur $(\Omega, \mathcal{F}, \mathbb{P})$. Alors les (X_i) sont indépendantes si et seulement si pour tout J fini, pour toutes fonctions bornées g_j telles que $g_j(X_j)$ soit intégrable, alors

$$\mathbb{E} \left[\prod_{j \in J} g_j(X_j) \right] = \prod_{j \in J} \mathbb{E} (g_j(X_j)).$$

Corollaire. (X_1, \dots, X_n) sont des variables aléatoires indépendantes si et seulement si pour tout $t \in \mathbb{R}^n$,

$$\varphi_{(X_1, \dots, X_n)}(t_1, \dots, t_n) = \prod_{j=1}^n \varphi_{X_j}(t_j).$$

2.5 Vecteurs aléatoires

Définition. On dit que X est un vecteur aléatoire sur $(\mathbb{A}, \mathcal{B}(\mathbb{R}^d))$, un espace probabilisé, si X est une fonction mesurable de \mathbb{A} dans $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$.

Définition. Soit X un vecteur aléatoire sur $(\mathbb{A}, \mathcal{B}(\mathbb{R}^d))$ à valeurs dans \mathbb{R}^d . Alors la loi (ou mesure) de probabilité de X , IP_X , est définie de façon unique à partir de la fonction de répartition de X tel que pour $x = (x_1, \dots, x_d)$,

$$F_X(x) = IP_X\left(\prod_{i=1}^d]-\infty, x_i]\right) = IP(X \in \prod_{i=1}^d]-\infty, x_i]).$$

Propriété. Soit X un vecteur aléatoire sur $(\mathbb{A}, \mathcal{B}(\mathbb{R}^d))$ à valeurs dans \mathbb{R}^d . On suppose que $X = (X_1, \dots, X_d)$. Alors les X_i sont des variables aléatoires sur $(\mathbb{A}, \mathcal{B}(\mathbb{R}^d))$, de fonction de répartition

$$F_{X_i}(x_i) = \lim_{x_j \rightarrow +\infty, j \neq i} F_X(x_1, \dots, x_i, \dots, x_d).$$

Les mesures de probabilité déterminées de façon unique à partir de x_i sont appelées lois marginales de X .

On se place maintenant dans la base canonique orthogonale $\{e_1, \dots, e_d\}$ de \mathbb{R}^d . On définit $IE(Z)$, le vecteur dont les coordonnées sont les espérances des coordonnées de Z . Ainsi, si dans la base canonique $Z = (Z_1, \dots, Z_d)'$,

$$IE(Z) = IE \begin{pmatrix} Z_1 \\ \vdots \\ Z_d \end{pmatrix} = \begin{pmatrix} IE(Z_1) \\ \vdots \\ IE(Z_d) \end{pmatrix}.$$

De la même manière, on définit l'espérance d'une matrice dont les coordonnées sont les covariances par la matrice dont les coordonnées sont les espérances de chacune de ces variables.

Ceci nous permet de définir la matrice de variance-covariance de Z de la manière suivante :

$$\text{var}(Z) = IE[(Z - IE(Z))(Z - IE(Z))']$$

donc si $Z = (Z_1, \dots, Z_d)'$,

$$\text{var}(Z) = \begin{pmatrix} \text{var}(Z_1) & \text{Cov}(Z_1, Z_2) & \dots & \text{Cov}(Z_1, Z_d) \\ \text{Cov}(Z_1, Z_2) & \text{var}(Z_2) & \dots & \text{Cov}(Z_2, Z_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(Z_1, Z_d) & \text{Cov}(Z_2, Z_d) & \dots & \text{var}(Z_d) \end{pmatrix}$$

matrice (d, d) dont les éléments diagonaux sont les variances et les éléments non diagonaux sont les covariances des coordonnées de Z (remarquons que la variance est égale à $\text{Cov}(Z_i, Z_i)$).

On vérifie également que si C est une matrice (p, d) à coefficients réels constants et si Z est un vecteur aléatoire à valeurs dans \mathbb{R}^d , alors CZ est un vecteur de taille p de matrice de variance-covariance

$$\text{var}(CZ) = C \cdot \text{var}(Z) \cdot C'.$$

En particulier, si p vaut 1, alors $C = h$ est un vecteur de taille d , et :

$$\text{var}(h \cdot Z) = h' \cdot \text{var}(Z) \cdot h.$$

Notons que cette dernière quantité est un scalaire. Soient Y_1, \dots, Y_d des variables aléatoires indépendantes de même loi $N(0, \sigma^2)$, indépendantes (ce qui, dans le cas gaussien, équivaut à $\text{Cov}(Y_i, Y_j) = 0$ pour $i \neq j$).

On considère le vecteur $Y_1 = (Y_1, \dots, Y_d)'$. En raison de l'indépendance, Y est un vecteur gaussien admettant

une densité f_Y (par rapport à la mesure de Lebesgue) qui satisfait le produit des densités de chacune des coordonnées, soit :

$$\begin{aligned} f_Y(y_1, \dots, y_d) &= f_{Y_1}(y_1) \times f_{Y_2}(y_2) \times \dots \times f_{Y_d}(y_d) \\ &= 2\pi\sigma^2^{-d/2} \exp -\frac{1}{2\sigma^2}(y_1^2 + \dots + y_d^2) \\ &= 2\pi\sigma^2^{-d/2} \exp -\frac{y^2}{2\sigma^2}, \end{aligned}$$

avec $y = (y_1, \dots, y_d)$. On voit donc que la densité ne dépend que de la norme euclidienne de y : elle est constante sur toutes les sphères centrées en 0. Cela implique qu'elle est invariante par rotation orthogonale d'axe passant par 0 : elle est invariante par toutes les rotations. On dira que Y suit une loi gaussienne isotrope. Rappelons que les coordonnées y_i sont indépendantes : à des changements de bases (BON). En conséquence, on a la propriété importante :

Propriété. Soit Y un vecteur aléatoire de loi normale isotrope variance σ^2 : c'est-à-dire que dans une BON les coordonnées de Y vérifient $E(Y) = 0$ et $\text{var}(Y) = \sigma^2 \cdot \text{Id}$. Alors les coordonnées de Y dans toute BON sont encore des $N(0, \sigma^2)$ indépendantes.

Voici maintenant l'un des résultats (encore à portée de Cochrane) que nous utiliserons le plus et nous en donnons donc une démonstration.

Théorème (Théorème de Cochran). Soient E_1 et E_2 , deux sous-espaces vectoriels orthogonaux de $E = \mathbb{R}^d$ de dimensions respectives k_1 et k_2 et soit Y un vecteur aléatoire de loi normale centrée isotrope de variance σ^2 . Alors $P_{E_1}(Y)$ et $P_{E_2}(Y)$ sont deux variables aléatoires gaussiennes indépendantes et $P_{E_1}(Y)^2$ (resp. $P_{E_2}(Y)^2$) est une loi $\chi^2(k_1)$ (resp. $\chi^2(k_2)$). Ce théorème se généralise naturellement pour 2 sous-espaces vectoriels orthogonaux de $E = \mathbb{R}^d$.

Démonstration. Soit (e_1, \dots, e_{k_1}) et $(e_{k_1+1}, \dots, e_{k_1+k_2})$ des bases orthonormales de E_1 et E_2 (respectivement). L'ensemble de ces deux bases étend une

$$(e_1, \dots, e_{k_1}, e_{k_1+1}, \dots, e_{k_1+k_2}, e_{k_1+k_2+1}, \dots, e_d)$$

pour former une BON de \mathbb{R}^d (du fait que E_1 et E_2 sont orthogonaux).

Soit (Y_1, \dots, Y_d) , les coordonnées de Y dans cette base. Elles sont indépendantes de $N(0, \sigma^2)$ car le changement de base est orthogonal et nous avons vu que la distribution de Y est isotrope. Comme

$$P_{E_1}(Y) = Y_1 e_1 + \dots + Y_{k_1} e_{k_1} \Rightarrow P_{E_1}(Y)^2 = \sigma^2 \left(\frac{Y_1^2}{\sigma^2} + \dots + \frac{Y_{k_1}^2}{\sigma^2} \right)$$

$$P_{E_2}(Y) = Y_{k_1+1} e_{k_1+1} + \dots + Y_{k_1+k_2} e_{k_1+k_2} \Rightarrow P_{E_2}(Y)^2 = \sigma^2 \left(\frac{Y_{k_1+1}^2}{\sigma^2} + \dots + \frac{Y_{k_1+k_2}^2}{\sigma^2} \right).$$

On voit bien ainsi l'indépendance entre les deux projections et que la loi de $P_{E_1}(Y)^2$ (resp. $P_{E_2}(Y)^2$) est une loi $\chi^2(k_1)$ (resp. $\chi^2(k_2)$). ■

On peut enfin plus généralement définir un vecteur gaussien Y à valeurs dans \mathbb{R}^d (c'est-à-dire, d'espérance $\mu \in \mathbb{R}^d$ et de matrice de variance-covariance Σ quelconques (du moment que Σ soit une matrice de Toeplitz définie positive)). Equivalently, on peut définir un vecteur aléatoire de densité par rapport à la mesure de Lebesgue sur \mathbb{R}^d ,

$$f_Y(y) = \frac{(2\pi)^{-n/2}}{|\Sigma|} \exp -\frac{1}{2}(y - \mu)' \cdot \Sigma^{-1} \cdot (y - \mu),$$

pour $y \in \mathbb{R}^d$, et avec Σ le déterminant de la matrice Σ . Remarquons une nouvelle fois que la variance-covariance est déterminée par la loi de probabilité d'un vecteur gaussien.

À partir des propriétés des vecteurs gaussiens, on obtient le fait que :

Propriété. Soit Y un vecteur gaussien à valeurs dans \mathbb{R}^d , d'espérance $\mu \in \mathbb{R}^d$ et de matrice de variance-covariance Σ . Soit C une matrice réelle $(p, d) \in \mathbb{R}^{p \times d}$. Alors $C \cdot Y$ est un vecteur gaussien tel que :

$$C \cdot Y \sim N(C \cdot \mu, C \cdot \Sigma \cdot C')$$

On en déduit les conséquences suivantes :

- si Y est un vecteur gaussien isotrope de variance σ^2 et h un vecteur de \mathbb{R}^d , alors $h \cdot Y$ est une combinaison linéaire des coordonnées de Y telle que :

$$h \cdot Y \text{ suit la loi } N(0, \sigma^2 \cdot h' \cdot h) = N(0, \sigma^2 \cdot \|h\|^2)$$

- si Y est un vecteur gaussien d'espérance μ et de matrice de variance Σ et si h un vecteur de \mathbb{R}^d , alors $h \cdot Y$ est une combinaison linéaire des coordonnées de Y telle que :

$$h \cdot Y \text{ suit la loi unidimensionnelle } N(h' \cdot \mu, h' \cdot \Sigma \cdot h)$$

(Pour une présentation plus détaillée des notions sur les vecteurs gaussiens on peut consulter le livre P. Toulouze, 1999, chapitre 2)

2.6 Fonctions caractéristiques

Définition. Soit X un vecteur réel à valeurs dans \mathbb{R}^d . La fonction caractéristique de X est la fonction $\varphi_X : \mathbb{R}^d \rightarrow \mathbb{C}$ telle que

$$\varphi_X(t) = \mathbb{E}[e^{i \langle t, X \rangle}] = \int_{\mathbb{R}^d} e^{i \langle t, x \rangle} d\mathbb{P}_X(x),$$

où $\langle \cdot, \cdot \rangle$ désigne le produit scalaire euclidien sur \mathbb{R}^d . Reliqu $\langle t, x \rangle = \sum_{i=1}^d t_i x_i$ pour $t = (t_1, \dots, t_d)$ et $x = (x_1, \dots, x_d)$.

Remarque. La fonction caractéristique existe sur \mathbb{R}^d et est aussi la transformée de Fourier de la mesure \mathbb{P}_X .

Théorème. Soit X et Y des vecteurs réels à valeurs dans \mathbb{R}^d , les lois \mathbb{P}_X et \mathbb{P}_Y . Alors $\mathbb{P}_X = \mathbb{P}_Y$ si et seulement si $\varphi_X = \varphi_Y$.

Théorème (Théorème d'inversion). Si X est un vecteur réel à valeurs dans \mathbb{R}^d et si φ_X est une fonction intégrable par rapport à la mesure de Lebesgue, alors X admet une densité par rapport à λ telle que pour $x \in \mathbb{R}^d$,

$$f_X(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{-i \langle t, x \rangle} \varphi_X(t) dt.$$

Proposition. Si X est une variable réelle sur $(\Omega, \mathcal{A}, \mathbb{P})$ de fonction caractéristique φ_X . Alors si $\mathbb{E}[|X|^n] < +\infty$ (ou $X \in L^n(\Omega, \mathcal{A}, \mathbb{P})$), φ_X est n -fois dérivable et $\varphi_X^{(n)}(t) = i^n \mathbb{E}[X^n e^{itX}]$.

Remarque. Lorsque ces moments existent, $\mathbb{E}[X^n] = \varphi_X^{(n)}(0)$.

2.7 Convergence de suites de variables réelles

Le même (Lemme de Borel-Carleson) ne s'applique pas aux événements sur $(\Omega, \mathcal{A}, \mathbb{P})$.

1. Si $\mathbb{P}(A_n) < +\infty$ alors $\mathbb{P}(\limsup A_n) = 0$.

2. Si les (A_n) sont indépendants, $\mathbb{P}(A_n) = +\infty$ implique que $\mathbb{P}(\limsup A_n) = 1$.

Définition. Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables réelles sur $(\Omega, \mathcal{A}, \mathbb{P})$. On dit que

(X_n) converge en probabilité vers X , noté $X_n \xrightarrow{P} X$, lorsque pour tout $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0.$$

\$(X_n)\$ converge dans \$(\mathbb{R}, \mathcal{A}, \mathbb{P})\$ vers \$X\$, si et seulement si \$\lim_{n \rightarrow +\infty} \mathbb{E}|X_n - X|^p = 0\$, avec \$p > 0\$, lors que

$$\lim_{n \rightarrow +\infty} \mathbb{E}|X_n - X|^p = 0.$$

\$(X_n)\$ converge en loi vers \$X\$, si et seulement si \$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x)\$ pour tout \$x \in \mathbb{R}\$ tel que \$F_X\$ continue en \$x\$.

$$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x) \text{ pour tout } x \in \mathbb{R} \text{ tel que } F_X \text{ continue en } x.$$

\$(X_n)\$ converge presque sûrement vers \$X\$, si et seulement si \$\mathbb{P}(\lim_{n \rightarrow +\infty} X_n = X) = 1\$ et que
 $\forall \omega \in E, \lim_{n \rightarrow +\infty} X_n(\omega) = X(\omega)$
 \Leftrightarrow pour tout \$\epsilon > 0\$,
 $\lim_{n \rightarrow +\infty} \mathbb{P}(\sup_{m \geq n} |X_m - X| > \epsilon) = 0.$

$$\lim_{n \rightarrow +\infty} \mathbb{P}(\sup_{m \geq n} |X_m - X| > \epsilon) = 0.$$

Propriété. 1. p.s. \$\mathbb{P} \rightarrow \mathbb{P} \rightarrow L\$.

2. pour \$p \geq 1\$, \$\mathbb{P} \rightarrow \mathbb{P} \rightarrow L^p\$.

3. La convergence en loi n'est pas la convergence en probabilité. \$(X_n) \xrightarrow{\mathbb{P}} C \Leftrightarrow (X_n) \xrightarrow{L} C\$ pour \$C\$ une constante.

4. Si \$g\$ est une fonction continue alors \$(X_n) \xrightarrow{\mathbb{P}} X \Rightarrow (g(X_n)) \xrightarrow{\mathbb{P}} g(X)\$.

Propriété. 1. Si pour tout \$\epsilon > 0\$, \$\sum_{n=0}^{\infty} \mathbb{P}(|X_n - X| > \epsilon) < +\infty\$ alors \$X_n \xrightarrow{\text{p.s.}} X\$ (application du Lemme de Borel-Cantelli).

2. Si il existe \$r > 0\$ tel que \$\sum_{n=0}^{\infty} \mathbb{E}|X_n - X|^r < +\infty\$ et \$\sum_{n=0}^{\infty} \mathbb{E}|X_n - X|^r < +\infty\$ alors \$X_n \xrightarrow{\text{p.s.}} X\$.

Théorème (Loi faible des Grands Nombres). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si \$\mathbb{E}|X_i| < +\infty\$,
 $\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\mathbb{P}} m = \mathbb{E}X_i.$

$$\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\mathbb{P}} m = \mathbb{E}X_i.$$

Théorème (Loi forte des Grands Nombres). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si \$\mathbb{E}|X_i| < +\infty\$,
 $\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\text{p.s.}} m = \mathbb{E}X_i.$

$$\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\text{p.s.}} m = \mathbb{E}X_i.$$

Théorème (Théorème de la limite centrale). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si \$\sigma^2 = \mathbb{E}X_i^2 < +\infty\$, et \$m = \mathbb{E}X_i\$,
 $\frac{\sqrt{n}(\overline{X}_n - m)}{\sigma} \xrightarrow{L} N(0, 1).$

$$\frac{\sqrt{n}(\overline{X}_n - m)}{\sigma} \xrightarrow{L} N(0, 1).$$

Théorème (Loi forte des Grands Nombres multidimensionnelle). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de vecteurs aléatoires à valeurs dans \$\mathbb{R}^d\$ indépendantes et identiquement distribuées. Alors si \$\mathbb{E}|X_i| < +\infty\$ (pour une norme sur \$\mathbb{R}^d\$),
 $\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\text{p.s.}} m = \mathbb{E}X_i.$

$$\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow{\text{p.s.}} m = \mathbb{E}X_i.$$

Théorème (Théorème de la limite centrale multidimensionnelle). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de vecteurs aléatoires à valeurs dans \$\mathbb{R}^d\$ indépendantes et identiquement distribuées. Alors si \$\Sigma\$ matrice de covariance de chaque \$X_i\$ existe, et \$m = \mathbb{E}X_i\$,
 $\sqrt{n}(\overline{X}_n - m) \xrightarrow{L} N_d(0, \Sigma).$

$$\sqrt{n}(\overline{X}_n - m) \xrightarrow{L} N_d(0, \Sigma).$$

Théorème (Delta-méthode). Soit \$(X_n)_{n \in \mathbb{N}}\$ une suite de vecteurs aléatoires à valeurs dans \$\mathbb{R}^d\$ indépendantes et identiquement distribuées, tel que \$\Sigma\$ matrice de covariance de chaque \$X_i\$ existe et \$m = \mathbb{E}X_i\$. Soit \$g: \mathbb{R}^d \rightarrow \mathbb{R}^p\$ une fonction de classe \$C^1\$ sur un voisinage autour de \$m\$, de matrice Jacobienne \$J_g(m)\$. Alors,
 $\sqrt{n}(g(\overline{X}_n) - g(m)) \xrightarrow{L} N_d(0, J_g(m) \cdot \Sigma \cdot J_g'(m)).$

$$\sqrt{n}(g(\overline{X}_n) - g(m)) \xrightarrow{L} N_d(0, J_g(m) \cdot \Sigma \cdot J_g'(m)).$$

2.8 Espérance conditionnelle

Définition. Soit Y une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$. \mathcal{B} est un sous-trièbre de \mathcal{A} . Soit $Y \in L^2(\Omega, \mathcal{A}, \mathbb{P})$. Alors on note $E(Y|\mathcal{B})$ la projection orthogonale de Y sur $L^2(\Omega, \mathcal{B}, \mathbb{P})$, appelée espérance conditionnelle de Y sachant \mathcal{B} . Ains :

$$E(Y|\mathcal{B}) = \inf_{Z \in L^2(\Omega, \mathcal{B}, \mathbb{P})} \|Y - Z\|_{L^2}^2.$$

Par extension, $Y \in L^1(\Omega, \mathcal{A}, \mathbb{P})$, on définit l'espérance conditionnelle par rapport à \mathcal{B} comme l'unique (p.s.) variable aléatoire B -mesurable vérifiant p.s. :

$$\int_B E(Y|\mathcal{B}) d\mathbb{P} = \int_B Y d\mathbb{P}, \quad \text{pour tout } B \in \mathcal{B}.$$

Définition. Par convention, si X un vecteur à valeurs dans \mathbb{R}^n sur $(\Omega, \mathcal{A}, \mathbb{P})$ et si Y une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, on note $E(Y|X) = E(Y|X^{-1}(\mathcal{B}(\mathbb{R}^n)))$.

Propriété. 1. Lemme de Doob : Pour $Y \in L^1(\Omega, \mathcal{A}, \mathbb{P})$, et X une v.a. de $(\Omega, \mathcal{A}, \mathbb{P})$, alors p.s. $E(Y|X) = h(X)$, avec h une fonction borélienne.

2. Pour Y_1 et Y_2 deux variables aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$, et $(a, b, c) \in \mathbb{R}^3$, alors

$$E(aY_1 + bY_2 + c|\mathcal{B}) = aE(Y_1|\mathcal{B}) + bE(Y_2|\mathcal{B}) + c.$$

3. Si $Y_1 \leq Y_2$, alors $E(Y_1|\mathcal{B}) \leq E(Y_2|\mathcal{B})$.

4. Le Lemme de Fatou et les théorèmes de Beppo-Levi et Lebesgue s'appliquent avec l'espérance conditionnelle.

5. Si $Y \in L^2(\Omega, \mathcal{B}, \mathbb{P})$, alors $E(Y|\mathcal{B}) = Y$; ains si $E(g(X)|X) = g(X)$ pour g une fonction mesurable réelle.

6. On a $E(E(Y|\mathcal{B})) = E(Y)$.

7. Si $Y^{-1}(\mathcal{B}(\mathbb{R}))$ et \mathcal{B} sont indépendants alors $E(Y|\mathcal{B}) = E(Y)$; ains si X et Y sont indépendants, $E(Y|X) = E(Y)$.

8. Si (X, Y) est un couple de v.a. à valeurs dans \mathbb{R}^2 possédant une densité $f_{(X,Y)}$ par rapport à la mesure de Lebesgue, alors si X est intégrable,

$$E(Y|X=x) = \frac{\int_{\mathbb{R}} y \cdot f_{(X,Y)}(x, y) dy}{\int_{\mathbb{R}} f_{(X,Y)}(x, y) dy}, \quad \text{pour tout } x \text{ tel que } \int_{\mathbb{R}} f_{(X,Y)}(x, y) dy > 0.$$

Proposition. Si (Y, X_1, \dots, X_n) est un vecteur gaussien, alors $E(Y, X_n) = a_0 + a_1 X_1 + \dots + a_n X_n$, où les a_i sont des réels.

3 Estimation par méthode

3.1 Définitions

Dans toute la suite, on se place sur $(\Omega, \mathcal{A}, \mathbb{P})$ un espace de probabilité et $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires, où chaque X_n est définie sur $(\Omega, \mathcal{A}, \mathbb{P})$ et est à valeurs dans \mathbb{R} .

Définition. - On appelle modèle statistique de dimension n un espace $\mathcal{A}'_n(\Omega)$, où \mathcal{A}'_n est une tribu sur $(\Omega, \mathcal{A}, \mathbb{P})$ et μ une mesure de probabilité sur $(\Omega^n, \mathcal{A}'_n)$.

- On appelle échantillon de taille n du modèle statistique (\mathcal{A}'_n, μ) le vecteur aléatoire (X_1, \dots, X_n) distribué selon la loi μ . Pour $\omega \in \Omega$, $(X_1(\omega), \dots, X_n(\omega))$ vecteur de \mathbb{R}^n est appelé échantillon observé. C'est à partir de ce vecteur qu'est tiré l'estimateur $\hat{\theta}$.

Définition. On appelle :

- Modèle statistique paramétré, une famille de modèles de la forme : $(\mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$.

- Modèle statistique semi-paramétré, une famille de modèles de la forme : $(\mathcal{A}'_n, \mathbb{P}_{\theta, f}, \theta \in \Theta, f \in F)$, où $\Theta \subset \mathbb{R}^p$ et F n'est pas de dimension finie.

-Modèle statistique non-paramétrique : une famille de mesures de la forme : $(\Omega, \mathcal{A}_n, \mathbb{P}_\theta, f \in F)$, où F n'est pas de dimension finie.

Définition. - On dit que le modèle paramétrique : $(\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, est dominé par une mesure μ lorsque est absolument continu par rapport à μ pour tout θ .
 - On se place dans le cadre d'un modèle paramétrique $(\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ pour $(x_1, \dots, x_n) \in (\Omega)^n$, la fonction $\theta \in \Theta \rightarrow L_\theta(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mu}(x_1, \dots, x_n)$ est appelée une vraisemblance déterministique.

Exemple. - Dans le cas où μ est la mesure de Lebesgue sur \mathbb{R} la vraisemblance sera la densité (classique) en (x_1, \dots, x_n) .

- Dans le cas où μ est comptable, la vraisemblance sera la probabilité (x_1, \dots, x_n) .

- Attention ! si le support dépend de θ , la mesure qui domine (ainsi qu'on ne peut dépendre de θ : il ne faut pas oublier de le préciser dans l'expression de la vraisemblance.

Définition. - Lorsque l'on dispose d'un échantillon (X_1, \dots, X_n) du modèle statistique $(\Omega, \mathcal{A}_n, \mu)$, une statistique est une application mesurable de $(\Omega)^n$ dans $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, donc un vecteur euclidien défini sur $(\Omega, \mathcal{A}_n, \mathbb{P})$ à valeur dans \mathbb{R}^d et tel que :

$$T_n = h(X_1, \dots, X_n), \quad \text{où } h: (\Omega)^n \rightarrow \mathbb{R}^d \text{ est mesurable.}$$

Exemple. - Estimateur du paramètre d'une loi de Bernoulli.

Estimateur de l'espérance et de la variance par la moyenne et la variance empirique.

Estimateurs du paramètre d'un échantillon (X_1, \dots, X_n) de loi uniforme sur $[0, \theta]$.

Test sur la moyenne.

3.2 Statistiques exhaustives

On se place désormais dans le cadre d'une famille statistique paramétrique $(\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominée par une mesure μ .

Exemple 1. Soit le modèle statistique paramétrique $[0, \infty[^n, \mathcal{B}([0, \infty[^n), \mathcal{U}([0, \theta]^n), \theta \in]0, +\infty[$. On dispose donc d'un échantillon (X_1, \dots, X_n) de v.a.i.i.d. suivant une loi uniforme sur $[0, \theta]$. Si on considère $\max\{X_1, \dots, X_n\}$ cela semble suffire pour caractériser l'information sur θ que contient (X_1, \dots, X_n) : on a donc une "information" sur θ contenue dans (X_1, \dots, X_n) , un vecteur de taille n , par une statistique de taille 1.

2. De même, si on considère le modèle statistique paramétrique $\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), \mathcal{B}(p)^{\otimes n}, p \in [0, 1]$ (on dispose donc d'un échantillon (X_1, \dots, X_n) de v.a.i.i.d. suivant une loi de Bernoulli de paramètre p) alors la statistique $X_1 + \dots + X_n$ contient toute l'"information" sur p contenue dans (X_1, \dots, X_n) .

Comment exprimer formellement ce fait qu'une statistique possède toute l'information sur le paramètre ?

Définition. Soit T une statistique déterministique paramétrique donnée par une valeur dans \mathbb{R}^d . On dit que T est une statistique exhaustive si elle est une statistique suffisante (donc dans $(\Omega)^n, \mathcal{A}_n, \mathbb{P}_\theta$) alors $\mathbb{E}_\theta(S | T)$ ne dépend (presque sûrement) pas de θ .

Théorème (Théorème de factorisation de Neyman). Soit (X_1, \dots, X_n) un n -échantillon et soit T une statistique déterministique paramétrique donnée par une valeur dans \mathbb{R}^d . Rudin $\in \mathbb{N}^*$. La statistique T est exhaustive si et seulement si il existe une fonction $\mathbb{R}^n \rightarrow \mathbb{R}_+$ et une fonction $g: \mathbb{R}^d \rightarrow \mathbb{R}_+$, tel que l'on puisse écrire pour tout $(x_1, \dots, x_n) \in (\Omega)^n$:

$$L_\theta(x_1, \dots, x_n) = g_\theta(T(x_1, \dots, x_n)) \cdot h(x_1, \dots, x_n) \quad \text{pour tout } \theta \in \Theta.$$

Le même soit le modèle statistique paramétrique $(\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$. Alors, ce modèle est dominé si et seulement si il existe une sous-famille dénombrable $(\mathbb{P}_i)_{i \in \mathbb{N}}$ tel que pour tout $\theta \in \Theta$, $\forall i \in \mathbb{N}$, $\mathbb{P}_\theta(A) = 0$ en tant que $\mathbb{P}_i(A) = 0$. Toute mesure de probabilité formée $\mathbb{P} = \sum_{i \in \mathbb{N}} \alpha_i \mathbb{P}_i$ avec $\alpha_i > 0$ pour tout $i \in \mathbb{N}$ et $\sum_{i \in \mathbb{N}} \alpha_i = 1$ domine le modèle.

Démonstration du lemme Il est bien clair que si une telle mesure existe, le lemme est démontré.
 \Rightarrow Montrons maintenant que si elle est donnée par une mesure μ alors la famille $(P_{\theta})_{\theta \in \Theta}$ est pexiste.
 En premier lieu, μ est une mesure non finie mais σ -finie (par exemple la mesure de Lebesgue),
 définie par $\mu(A) = \sum_{i=1}^{\infty} \frac{1}{2^i} \mu(A \cap A_i)$ pour tout $A \in \mathcal{A}$, est une mesure de probabilité à l'égard de μ (avec $(A_i)_{i \in \mathbb{N}^*}$ une partition de Ω telle que $0 < \mu(A_i) < \infty$ pour tout $i \in \mathbb{N}^*$). On travaille donc désormais avec μ .
 Pour $\theta \in \Theta$, soit B_{θ} le sous-ensemble de Ω^n qui est le support de la densité P_{θ} par rapport à μ .
 Soit

$$C = \bigcup_{i \in I} B_{\theta_i}, \quad I \subset \mathbb{N}, \theta_i \in \Theta,$$

l'ensemble de toutes les unions finies d'ensembles B_{θ} . Soit $M = \sup_{C \in \mathcal{C}} \mu_P(C)$. Soit $(G_n)_{n \in \mathbb{N}}$ une suite d'ensembles C telle que la suite $(\mu(G_n))_n$ converge vers M (une telle suite existe forcément si M n'est pas le supremum). Remarquons que $\mu_P(G_n) \leq M$ et donc $\mu_P(G_n) \rightarrow M$.
 Soit $(\theta_n)_{n \in \mathbb{N}}$ une suite de Θ telle que $G_n = \bigcup_{k=1}^n B_{\theta_k}$. Si on pose :

$$D = \bigcup_{n \in \mathbb{N}} G_n = \bigcup_{k \in \mathbb{N}} B_{\theta_k},$$

alors $M = \mu_P(D)$ et pour tout $\theta \in \Theta, B_{\theta} \cup D \in \mathcal{C}$ et :

$$\mu_P(B_{\theta} \cup D) \leq M \leq \mu_P(B_{\theta} \cup D) = \mu_P(B_{\theta} \cap D^c) + \mu_P(D)$$

Donc pour tout $\theta \in \Theta, \mu_P(B_{\theta} \cap D^c) = 0$ so $\forall \theta \in \Theta, \mu_P(B_{\theta} \cap D^c) = 0$ puis que $\mu_P < \mu$. En conséquence, pour tout $A \in \mathcal{A}_n, A \subset B_{\theta} \cup D^c = (\Omega')^n$, soit :

$$\mu_P(A \cap D^c) = 0, \quad \text{car par définition des } B_{\theta} \quad \mu_P(B_{\theta} \cap D^c) = 0.$$

Si on suppose maintenant que $\mu_P(A) = 0$, avec la suite précédemment définie, alors $\mu_P(A \cap B_{\theta_k}) = 0$ par définition des B_{θ} et donc $\mu_P(A \cap D) = 0$ (par la propriété de σ -additivité d'une mesure). Comme $\mu_P < \mu$, on en déduit que $\forall \theta \in \Theta, \mu_P(A \cap D) = 0$ et donc $\mu_P(A) = \mu_P(A \cap D) + \mu_P(A \cap D^c) = 0$.
 Ainsi, μ_P domine bien μ pour tout $\theta \in \Theta$. ■

Démonstration du théorème de factorisation de Neyman : Soit $(P_{\theta})_{\theta \in \Theta}$ une mesure de probabilité dominante construite comme dans le lemme.

\Leftarrow Soit $(T(x)) \cdot h(x)$ avec $C \in (\Omega')^n$ est la densité P par rapport à μ , alors $(\sum_{i=1}^n a_i \cdot g_{\theta_i}(T(x)) \cdot h(x) = g_*(T(x)) \cdot h(x)$ est une densité P^* par rapport à μ . Alors, comme $(T(x)) \cdot h(x) > 0$ p.s., donc p.s., pour toute variable aléatoire S telle que :

$$\begin{aligned} \mathbb{E}_{\theta}(S \cdot \mathbb{I}_B) &= \int_B S d\mu_P, \quad \text{pour tout } B \in \sigma(T), \text{ tribu engendrée par } T \\ &= \int_B S(x) \cdot g_{\theta}(T(x)) \cdot h(x) d\mu(x) \\ &= \int_B S(x) \cdot \frac{g_{\theta}(T(x)) \cdot h(x)}{g_*(T(x)) \cdot h(x)} d\mu_P^*(x) \\ &= \mathbb{E}_{P^*} \left(S \cdot \frac{g_{\theta}(T)}{g_*(T)} \right) \\ &= \mathbb{E}_{P^*} \left(S \cdot \frac{g_{\theta}(T)}{g_*(T)} \right) \cdot \mathbb{E}_{P^*}(S | T) \quad (\text{d'après la définition de l'espérance conditionnelle}) \\ &= \mathbb{E}_{P^*} \left(S \cdot \frac{g_{\theta}(T)}{g_*(T)} \right) \cdot \mathbb{E}_{P^*}(S | T) \end{aligned}$$

En conséquence, d'après la définition de l'espérance conditionnelle de $(S | T)$ par rapport à μ_P , on a $\mathbb{E}_{P^*}(S | T) = \mathbb{E}_{\theta}(S | T)$: la statistique T est bien exhaustive.

\Rightarrow On suppose que T est une statistique exhaustive pour la famille $(P_{\theta})_{\theta \in \Theta}$ et pour toute statistique S telle que :

$S, \forall \theta, \mathbb{E}(S | T) = \mathbb{E}_\theta(S | T)$. En conséquence, si on note $\varphi(x, \theta) = \frac{d\mathbb{P}_\theta}{d\mathbb{P}_*}(x)$ la densité de \mathbb{P}_θ par rapport à \mathbb{P}_* ,

$$\begin{aligned} \mathbb{E}_\theta(S) &= \mathbb{E}_\theta(\mathbb{E}_*(S | T)) \quad (\text{car } T \text{ est exhaustive et les propriétés de l'espérance conditionnelle}) \\ &= \mathbb{E}_* \int \varphi(X, \theta) \mathbb{E}_*(S | T) d\mathbb{P}_*, \quad \text{où } X \sim \mathbb{P}_* \\ &= \mathbb{E}_* \int \varphi(X, \theta) \mathbb{E}_*(S | T) dT, \quad (\text{d'après les propriétés de l'espérance conditionnelle}) \\ &= \mathbb{E}_* \int \varphi(X, \theta) T d\mathbb{E}_*(S | T), \quad (\text{car } \mathbb{E}_*(S | T) \text{ est une fonction de } T) \\ &= \mathbb{E}_* \int \varphi(X, \theta) T dT \\ &= \mathbb{E}_* \int S \varphi(X, \theta) dT \end{aligned}$$

Ainsi, la variable $\mathbb{E}_*(S | T)$ est une fonction de T (qui est elle-même une fonction sur $(\Omega')^n$), est la densité \mathbb{P}_θ par rapport à \mathbb{P}_* (par suite) la vraisemblance, est la densité \mathbb{P}_θ par rapport à \mathbb{P}_* , c'est-à-dire :

$$L_\theta(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mu}(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mathbb{P}_*}(x_1, \dots, x_n) \cdot \frac{d\mathbb{P}_*}{d\mu}(x_1, \dots, x_n) = \mathbb{E}_* \varphi(X, \theta) T \cdot h(x_1, \dots, x_n),$$

avec h une fonction mesurable. ■

Exemple. Différentes statistiques exhaustives pour les paramètres de loi normale, loi de Bernoulli, de loi gaussienne...

Propriété. On se place dans le cadre d'un espace statistique donné :

1. La statistique $T = T(X, \dots, X)$ est exhaustive.
2. Si T est une statistique exhaustive, existe-t-il une fonctionnelle h telle qu'une autre statistique U vérifie $T = h(U)$, alors U est également exhaustive.

On vient de voir que l'on peut toujours trouver une statistique exhaustive (même par exemple). Comme on a un théorème à vouloir le "maximum d'information" dans une statistique exhaustive, le problème est de savoir laquelle dimension minimale peut avoir cette statistique. En particulier, si $d = 1$, peut-on toujours trouver une statistique exhaustive de taille 1 ? L'exemple suivant montre que ce n'est pas toujours le cas :

Exemple. Soit le modèle statistique $\mathcal{E} = \{f_\theta : [0, \infty[\rightarrow \mathbb{R}, (\mathbb{P}_\theta)^{\otimes n}, \theta \in \mathbb{R}_+\}$, où la densité \mathbb{P}_θ par rapport à la mesure de Lebesgue est $f_\theta(x) = \theta(e^\theta - 1)^{-1} e^{-\theta \cdot x}$, $x \in [0, \theta]$. Alors les statistiques $T_1 = X_1, \dots, X_n$ et $T_2 = X_1 + \dots + X_n$ ne sont pas chacune exhaustive alors que T_2 est exhaustive. On pourra même montrer que cette statistique est de taille minimale...

Définition. Une statistique exhaustive T est dite minimale si pour toute autre statistique exhaustive U est telle qu'il existe une fonctionnelle h vérifiant :

$$T = h(U).$$

Proposition. Soit un modèle statistique \mathcal{E} et soit $L_\theta(x_1, \dots, x_n)$ sa vraisemblance. Alors T est une statistique exhaustive minimale pour \mathcal{E} si et seulement si $L_\theta(x_1, \dots, x_n) = L_\theta(y_1, \dots, y_n)$ pour tout $(x_1, \dots, x_n) \in (\Omega')^n$ et $(y_1, \dots, y_n) \in (\Omega')^n$ tels que $T(x) = T(y)$.

$$\theta \mapsto \frac{L_\theta(x_1, \dots, x_n)}{L_\theta(y_1, \dots, y_n)} \text{ ne dépend pas de } \theta \iff T(x_1, \dots, x_n) = T(y_1, \dots, y_n). \quad (2)$$

Démonstration de la proposition. On suppose que (2) est vraie et on suppose (sans perte de généralité) que la vraisemblance est strictement positive sur $(\Omega')^n$. Notons $x^{(t)} \in T^{-1}(\{t\}) \subset (\Omega')^n$. Alors $\forall x \in T^{-1}(\{t\})$, $T(x) = T(x^{(t)})$ et donc d'après (2),

$$h(x) = \frac{L_\theta(x)}{L_\theta(x^{(t)})} \text{ est indépendant de } \theta.$$

Pour tout $x \in \Omega$, $L_\theta(x) = L_\theta(x^{(t)})$. Alors $L_\theta(x) = g_\theta(T(x))h(x)$. Comme ceci est vrai pour tout $x \in \Omega$, la statistique T est bien exhaustive.

Supposons maintenant que S est une autre statistique exhaustive. Alors, par le théorème de factorisation de Neyman, il existe deux fonctions $g_\theta^{(s)}$ et $h^{(s)}$ (ne dépendant pas de θ) telles que pour tout $x \in \Omega$, $L_\theta(x) = g_\theta^{(s)}(S(x)) \cdot h^{(s)}(x)$. Ainsi pour tout $x, y \in \Omega^n$ tels que $S(x) = S(y)$, alors :

$$\frac{L_\theta(x)}{L_\theta(y)} = \frac{g_\theta^{(s)}(S(x)) \cdot h^{(s)}(x)}{g_\theta^{(s)}(S(y)) \cdot h^{(s)}(y)} = \frac{h^{(s)}(x)}{h^{(s)}(y)}, \text{ qui est indépendant de } \theta.$$

Mais d'après (2) ceci n'est possible que si $T(x) = T(y)$. Donc T est une fonction de S et la statistique T est donc minimale. ■

Quelle serait une sorte d'opposé de la notion de statistique exhaustive minimale ? Cette notion de statistique ne dépendant pas du paramètre :

Définition Une statistique T d'un paramètre θ est dite libre si sa loi ne dépend pas du paramètre.

Or, de façon assez surprenante il peut arriver qu'une statistique exhaustive minimale coïncide avec une statistique libre qui intuitivement ne devrait pas être prise en compte pour donner toute l'information sur θ (soit par exemple la loi d'un échantillon de taille 2, la statistique $(X_1 - X_2, X_1 + X_2)$ est exhaustive minimale, mais $X_1 + X_2$ est libre). Aussi peut-on rajouter une autre caractérisation des statistiques exhaustives pour pouvoir atteindre une forme d'optimalité : ces statistiques satisfont qu'une fonctionnelle non constante de la statistique libre peut être augmentée sans modifier la loi.

Définition Une statistique exhaustive T est dite libre si pour toute fonctionnelle $h : \mathbb{R} \rightarrow \mathbb{R}$ telle que $h(T)$ soit intégrable, alors :

$$\forall \theta \in \Theta, \mathbb{E}_\theta(h(T)) = 0 \implies h(T) = 0.$$

Propriété. Soit un modèle statistique paramétrique donné :

1. si T est une statistique exhaustive et si pour toute fonctionnelle h bijective $h(T)$ est une statistique exhaustive compl.
2. si T est une statistique exhaustive et si T est un(e) statistique exhaustive minimale.
3. (Théorème de Basu) Si T est une statistique exhaustive et si T est libre sur le modèle.

Démonstration de la propriété 3. Théorème de Basu. Soit S une statistique libre pour le modèle $(\Omega, \mathcal{F}, \mathbb{P}, \theta \in \Theta \subset \mathbb{R}^p)$ et T une fonctionnelle telle que $(\mathbb{E}_\theta(h(T)))$ existe. Comme S est libre, on peut noter $e(f) = \mathbb{E}_\mu(f(S))$ une application linéaire ne dépendant pas de θ . Par suite, la statistique $(f(S) | T) - e(f)$ est une fonction de T mesurable telle que $\mathbb{E}_\mu(f(S) | T) - e(f) = 0$ pour tout $\theta \in \Theta$. Comme on a supposé que T est exhaustive compl., alors $\mathbb{E}_\theta(f(S) | T) = e(f)$ presque-sûrement : les statistiques S et T sont indépendantes. ■

Définition On suppose un paramétrique $(\Omega, \mathcal{A}, \mathbb{P}, \theta \in \Theta \subset \mathbb{R}^p)$ donné par un ensemble μ . Si, pour tout $(x_1, \dots, x_n) \in \Omega^n$ et $\theta \in \Theta$, la vraisemblance de ce vecteur par rapport à μ peut s'écrire sous la forme :

$$L_\theta(x_1, \dots, x_n) = \exp \left[\beta(\theta) + b(x_1, \dots, x_n) + \sum_{j=1}^p \alpha_j(x_1, \dots, x_n) \cdot \alpha_j(\theta) \right], \quad (3)$$

avec les fonctions $a : \Omega^n \rightarrow \mathbb{R}$, $b : \Omega^n \rightarrow \mathbb{R}$, $\alpha_j : \Theta \subset \mathbb{R}^p \rightarrow \mathbb{R}$, et $\beta : \Theta \rightarrow \mathbb{R}$, alors on dit que le modèle est exponentiel (quel qu'il appartient à la famille exponentielle).

Exemples Appartiennent à la famille exponentielle les lois :

- Loi discrètes : Lois de Bernoulli, binomiales, de Poisson, ...
- Loi "continues" : Lois normales, exponentielles, gamma, du chi-deux, ...

Remarque. Si (X_1, \dots, X_n) est un échantillon d'un modèle exponentiel (avec θ dans Θ) alors l'ensemble des valeurs prises par (X_1, \dots, X_n) ne dépend pas du paramètre θ .

Propriété. Soit un modèle exponentiel. Si pour $\theta \in \Theta$ on note $\alpha(\theta) = (a_1(\theta), \dots, a_p(\theta))$ et si l'ensemble $\alpha(\Theta)$ est d'intérieur non vide, alors $T(x, x_n) = (a_1(x_1, \dots, x_n), \dots, a_p(x_1, \dots, x_n))$ est une statistique exhaustive minimale et complète.

Démonstration de la propriété. Soit $g : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $E_\theta(g(T)) = 0$. On a $\theta \in \Theta$,

$$E_\theta(g(T)) = \int_{(\Omega')^n} g(T(x)) \cdot \exp(\beta(\theta) + b(x) + \langle T(x), \alpha(\theta) \rangle) d\mu(x),$$

où $\langle \cdot, \cdot \rangle$ désigne le produit scalaire. En considérant la mesure ν de densité $\exp(b(x))$ par rapport à μ , on obtient :

$$\begin{aligned} E_\theta(g(T)) = 0 &\Rightarrow \int_{(\Omega')^n} g(T(x)) \cdot \exp(\langle T(x), \alpha(\theta) \rangle) d\nu(x) = 0 \\ &\Rightarrow \int_{T((\Omega')^n)} g(y) \exp(\langle y, \alpha(\theta) \rangle) d\nu(y) = 0 \end{aligned}$$

pour tout $\theta \in \Theta$, en ayant noté la mesure image de ν par T et avec $T \in \mathbb{R}^p$. Si on note γ^+ et γ^- les parties positives et négatives de g (donc $g = \gamma^+ - \gamma^-$), et π^+ et π^- les mesures de densité γ^+ par rapport à ν alors, pour tout $\theta \in \Theta$:

$$\int_{T((\Omega')^n)} \exp(\langle y, \alpha(\theta) \rangle) d\pi^+(y) = \int_{T((\Omega')^n)} \exp(\langle y, \alpha(\theta) \rangle) d\pi^-(y).$$

En conséquence sur Θ , donc sur une partie d'intérieur non vide, les mesures π^+ ont des trajectoires communes. La propriété est : ces deux mesures sont égales et donc $\gamma^+ = \gamma^-$ presque partout (ce qui revient à $g = 0$). À partir des expressions des mesures, on montre que $g = 0$ presque partout. ■

3.3 Information de Fisher

Pour mesurer l'information fournie par un échantillon d'un modèle (ou une statistique sur ce modèle) au sujet d'un paramètre, une façon naturelle sera de mesurer comment varie localement la mesure de probabilité. On va encore savoir à l'aide de la notion de dérivée des moyennes de cette variable aléatoire. Pour ce faire on considère, lorsqu'il existe (pour x_1, \dots, x_n), et on s'intéressera à la matrice de covariance $(L_{ij}(x_1, \dots, x_n))$, dont on peut montrer qu'elle ne dépend pas du choix de la mesure dominante choisie. On introduit d'abord la notion de dérivée qui nous permettra d'exprimer cette quantité d'information.

Défini ti o n. Dans le cadre d'un modèle statistique paramétré $(\Omega, \mathcal{A}, \mathbb{P}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, donné par une mesure μ , on dira que θ est dérivable lorsque :

1. Θ est un ouvert de \mathbb{R}^p
2. la vraisemblance $L_\theta(x_1, \dots, x_n) \in (\Omega')^n, \forall \theta \in \Theta, L_\theta(x_1, \dots, x_n) > 0$;
3. $\forall (x_1, \dots, x_n) \in (\Omega')^n$, la fonction $\theta \in \Theta \rightarrow \log(L_\theta(x_1, \dots, x_n))$ est différentiable sur Θ par rapport à θ , et son gradient appartient à $\mathbb{R}^p, \forall \theta \in \Theta$;
4. $\forall \theta \in \Theta$, pour toute fonction $f : \mathbb{R}^p \rightarrow \mathbb{R}$ appartenant à $L^1((\Omega')^n, \mathcal{A}, \mathbb{P})$, alors :

$$\frac{\partial}{\partial \theta} \int_{(\Omega')^n} h(x) \cdot L_\theta(x) d\mu(x) = \int_{(\Omega')^n} h(x) \cdot \frac{\partial}{\partial \theta} L_\theta(x) d\mu(x). \quad (4)$$

Conséquence. Pour un modèle régulier, $E_{\theta_0}(\text{grad}(\log L_{\theta_0}(\cdot))) = 0$.

Démonstration. On a $E_{\mu}(L_\theta(\cdot)) = 1$ donc $E_{\mu}(\text{grad}(\log L_\theta(\cdot))) = 0$. Par conséquent $E_{\mu} \frac{\text{grad}(L_\theta(\cdot))}{L_\theta(\cdot)} = 0$, soit $E_{\theta_0}(\text{grad}(\log L_{\theta_0}(\cdot))) = 0$. ■

Défini ti o u n m o d è s t a t i s t i q u e p a r a q u e d o n n é e i r r é g u l i e r, o n a p p e l l e i n f o r m a t i o n d e F i s h e r, l a m a t r i c e :

$$I_n(\theta) = \mathbb{E}_{\theta} \sum_{1 \leq i, j \leq n} \frac{\partial (\log \text{gl}(X_1, \dots, X_n))}{\partial \theta_i} \times \frac{\partial (\log \text{gl}(X_1, \dots, X_n))}{\partial \theta_j}.$$

Pro pri  t  . Pou r un mod  le statistique param  trique donn  , et $\forall (x_1, \dots, x_n) \in (\Omega)^n$, la fonction $\theta \in \Theta \rightarrow \log(\mathcal{L}(\cdot))$ est $\mathcal{C}^2(\Theta)$, alors:

$$I_n(\theta) = -E_{\theta} \frac{\partial^2 (\log L_{\theta}(X_1, \dots, X_n))}{\partial \theta_i \cdot \partial \theta_j} \quad 1 \leq i, j \leq p.$$

Défini ti o n l' i n f o r m a t i o n d e F i s h e r $I(\theta)$ a s s o c i é e ` a u n e s t a t i s t i q u e T , l e x i s t e n c e d e l a m a t r i c e d e F i s h e r d e l a v r a i s e m b l a n c e d e T e s t e q u i v a l e n t ` a p a r t i r d e l a v r a i s e m b l a n c e d e T).

Propriété. Pour un module régulier, T est un étatistique libre si et seulement si $T(\phi) = 0$.

Démonstration : Si T est libre alors sa loi dépend pas de θ donc le gradient du log-likelihood de sa vraisemblance est nul ; l'information de Fisher est nulle.

\Leftarrow Si $Sil_n(\theta) = 0$, donc la statistique $g_n(\log_l(T))$ est centrée de variance nulle. Ainsi, pour tout $\theta \in \Theta$, il existe un ensemble N_θ de mesure 1 pour la mesure de probabilité μ_θ à T (donc, d'après la première hypothèse d'un modèle régulier tel que $\mu(N) = 1$) et tel que pour tout $t \in N_\theta$, $g_n(\log_l(t)) = 0$. Pour montrer que $g_n(\log_l(t)) = 0$ est bien une variable aléatoire μ -p.s., et donc que log_l est une fonction constante μ -p.s., nous faut montrer que finalement les θ_n dépendent pas de θ . Soit $\theta^{(d)} = \{\theta_i^{(d)}\}_{i \in \mathbb{N}}$ un sous-ensemble dénombrable dense dans Θ . Comme $\theta^{(d)}$ est dénombrable, il est clair que $N = \bigcap_{i \in \mathbb{N}} N_{\theta_i^{(d)}}$ est tel que $\mu(N) = 1$. De plus, pour tout $\theta \in \Theta$, il existe une sous-suite $(\theta_{n_k})_{k \in \mathbb{N}}$ de $\theta^{(d)}$ convergeant vers θ et telle que pour tout $i \in \mathbb{N}$, $g_n(\log_l(t)) = 0$. Comme une telle fonction g_n est continue, cette propriété passe à la limite, et donc pour tout $i \in \mathbb{N}$, $\forall \theta \in \Theta$, $g_n(\log_l(t)) = 0$. Comme N ne dépend pas de θ , la fonction $\theta \mapsto log_l(\cdot)$ est une constante μ -p.s. : la statistique T est bien libre. ■

Propriété. Pour un modèle régulier, si T est une statistique exhaustive $T(\theta) \in \mathcal{T}$ pour tout $\theta \in \Theta$.

Démonstration Comme T est une statistique exhaustive, peut écrie d'après la démonstration du Théorème de factorisation de Neyman que pour tout $t \in \mathcal{K}^{(n)}$ et tout $\theta \in \Theta$:

$$\frac{dP_{\theta}}{dP^*}(x_1, \dots, x) = g_{\theta}(T(x_1, \dots, x)).$$

On peut écrire cela pour la classe \mathbb{T} sous la forme: $\frac{dip_0^T}{dip^{*T}}(t) = g(t)$, pour tout $t \in T((\Omega')^n)$ et tout $\theta \in \Theta$. En conséquence, pour tout $t \in \mathbb{T}$

$$\begin{aligned}
I(\theta) &= \mathbb{E}_{\theta} \frac{\partial (\log g_{\theta}(X_1, \dots, X_N))}{\partial \theta_i} \times \frac{\partial (\log g_{\theta}(X_1, \dots, X_N))}{\partial \theta_j} \quad 1 \leq i, j \leq p \\
&= \frac{\partial (\log g_{\theta}(x))}{\partial \theta_i} \times \frac{\partial (\log g_{\theta}(x))}{\partial \theta_j} dP_{\theta}(x) \quad 1 \leq i, j \leq p \\
&= \frac{\partial (\log q_{\theta}(T(x)))}{\partial \theta_i} \times \frac{\partial (\log q_{\theta}(T(x)))}{\partial \theta_j} q_{\theta}(T(x)) dP^*(x) \quad \text{car } \log g_{\theta}(x) = \log q_{\theta}(T(x)) + \log h(x) \\
&= \frac{\partial (\log q_{\theta}(t))}{\partial \theta_i} \times \frac{\partial (\log q_{\theta}(t))}{\partial \theta_j} q_{\theta}(t) dP^{*T}(x) \quad \text{d'après le théorème du transport} \\
&= \frac{\partial (\log q_{\theta}(t))}{\partial \theta_i} \times \frac{\partial (\log q_{\theta}(t))}{\partial \theta_j} dP_{\theta}(t) \quad 1 \leq i, j \leq p \\
&= I_n^T(\theta). \quad \blacksquare
\end{aligned}$$

Remarque. En rajoutant certaines propriétés de T , on peut montrer que la réciproque est également vraie, et donc $c_T(\theta) = 0$ si et seulement si la statistique T est exhaustive.

Ainsi, on retrouve l'idée de la notation d'information de Fisher les "informations" qu'on peut tirer de la section de données. Voyons maintenant les applications de la notation de l'estimation par la méthode.

3.4 Application à l'estimation paramétrique

On se place dans le cadre d'un modèle statistique paramétrique $((\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$, défini par une mesure μ . Par ailleurs, on suppose que Θ est un ouvert.

Définition. – Soit $g: \Theta \rightarrow \mathbb{R}^p$, où $\theta \in \mathbb{R}^p$ avec $p \in \mathbb{N}^*$, une fonction mesurable. On appelle estimateur de la fonction g du paramètre $g(\theta)$ une statistique T à valeurs dans \mathbb{R}^p par laquelle on estime $g(\theta)$ est une statistique à valeurs dans \mathbb{R}^p est une réalisation de T .

- On appelle biais d'un estimateur T de $g(\theta)$ le vecteur constant $B(\theta) = \mathbb{E}(T) - g(\theta)$. On dira que l'estimateur est sans biais si $B(\theta) = 0$ pour tout θ .
- On appelle risque quadratique de l'estimateur T de $g(\theta)$ si $R(\theta) = \mathbb{E}(\|T - g(\theta)\|^2)$, où $\|\cdot\|$ désigne usuellement la norme euclidienne (mais peut être autre fonctionnelle positive et convexe). Si l'estimateur est sans biais alors, $R(\theta) = \text{Trace}(\text{cov}(T))$.

Pour pouvoir parler du comportement asymptotique d'une statistique, on va devoir se placer dans un modèle dans lequel on a une suite de variables aléatoires quel que soit le modèle pour lequel on se place $((\Omega)^{\mathbb{N}}, \mathcal{A}_{\mathbb{N}}, \mathbb{P}_\theta^{\mathbb{N}}, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$ (la dimension du paramètre reste constante pour un n fixé, une statistique sera d'abord une projection du "gros ensemble" de la loi puis une statistique "normale". On devra donc parler d'une suite d'estimateurs (T_n)

Définition. Pour un modèle statistique paramétrique $((\Omega^{\mathbb{N}}, \mathcal{A}_{\mathbb{N}}, \mathbb{P}_\theta^{\mathbb{N}}, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$, et pour (T_n) une suite d'estimateurs de $g(\theta)$:

- Si $\lim_{n \rightarrow \infty} B_n(\theta) = 0$, on dit que l'estimateur est asymptotiquement sans biais.
- On dit que (T_n) est convergent lorsque $T_n \xrightarrow{P} g(\theta)$.
- S'il existe une suite de réels positifs tels que $(T_n - g(\theta)) \xrightarrow[n \rightarrow +\infty]{L} Z_\theta$, où Z_θ est une loi centrée non nulle (ne dépend pas de n), on dit que T_n converge vers $g(\theta)$ à la vitesse

A priori être sans biais n'est pas un bon critère pour garantir une certaine efficacité d'un estimateur. On préfère pouvoir discriminer entre de potentiels estimateurs à l'aide d'un critère de risque quadratique ou sur la matrice de variance-covariance. Cependant, il n'y a pas de raison pour choisir un "meilleur" estimateur en ce sens. Pour en obtenir, on devra se limiter à une certaine classe d'estimateurs, celle des estimateurs sans biais.

Définition. Soit un modèle statistique paramétrique $((\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, et soit T un estimateur sans biais de $g(\theta)$. On dit que T est de variance minimum parmi les estimateurs sans biais de $g(\theta)$ lorsque pour tout estimateur sans biais de $g(\theta)$ on a $\text{cov}(T) \leq \text{cov}(S)$ (au sens où $\text{cov}(T) - \text{cov}(S)$ est une matrice positive).

Propriété. Si T est un estimateur de variance minimum parmi les estimateurs sans biais, alors il est unique à \mathbb{P} -p.s.

Démonstration. Soit S un autre estimateur que l'on suppose également de variance minimum parmi les estimateurs sans biais. Montrons d'abord que $\mathbb{E}(T - S) = 0$. En effet, si T est de variance minimum, en utilisant les matrices symétriques:

$$\begin{aligned} \text{cov}(T) &\leq \text{cov}(T + \alpha(T - S)) \\ &\leq \text{cov}(T) + 2\alpha \text{cov}(T - S) + \alpha^2 \text{cov}(T - S) \\ \Rightarrow 0 &\leq \alpha \cdot \alpha \cdot \text{cov}(S) + 2\alpha \mathbb{E}((T - S)) \quad \text{pour tout } \alpha \in \mathbb{R}. \end{aligned}$$

Comme $\text{cov}(T, S)$ est une matrice positive, on a $\text{cov}(T, S) = 0$ si et seulement si $T = S$ sur un ensemble de probabilité 1. Par suite, comme $\text{cov}(T, S) = 0$, on a $T = S$ sur un ensemble de probabilité 1.

Théorème (Rao-Blackwell) Soit T un estimateur sans biais de $g(\theta)$ et si S est une statistique exhaustive, alors $R = E(T | S)$, qui ne dépend pas de θ car S est exhaustive, est un estimateur sans biais de $g(\theta)$ de matrice de covariance inférieure ou égale à celle de T .

Démonstration Il est clair que $E(R) = E(T) = g(\theta)$. De plus, pour tout $u \in \mathbb{R}^p$ (avec $g: \mathbb{R} \rightarrow \mathbb{R}^p$),

$$\begin{aligned} \text{cov}(u, T) &= E_{\theta} (u \cdot (T - g(\theta))) \\ &= E_{\theta} (E_{\theta} (u \cdot (T - g(\theta)) | S)) \\ &\geq E_{\theta} (E_{\theta} (u \cdot (T - g(\theta)) | S))^2 \quad \text{d'après l'identité de Jensen,} \\ &\geq \text{cov}(u, R). \end{aligned}$$

Cela revient bien à dire que $\text{cov}(T) \geq \text{cov}(R)$.

Théorème (Lehmann-Scheffé) Si T est un estimateur sans biais de $g(\theta)$ et si S est une statistique exhaustive et complète, alors l'unique estimateur de $g(\theta)$ sans biais et à variance minimale est $R = E_{\theta}(T | S)$ (c'est-à-dire que R est une fonction de S).

Démonstration Soit T' un autre estimateur sans biais de $g(\theta)$. Soit $R = E(T' | S)$, on sait que $\text{cov}(T) \geq \text{cov}(R)$ d'après le théorème de Rao-Blackwell. Or $E(T' - R) = 0$ pour tout θ car les deux estimateurs sont sans biais. De plus comme R est une fonction de S , R est aussi une fonction de T' , et du fait que S est une statistique exhaustive et complète, on a $R = T'$ p.s. Par conséquent, pour tout θ , $\text{cov}(R) = \text{cov}(T')$ et donc $\text{cov}(R) = \text{cov}(T)$: R est bien l'estimateur sans biais de variance minimale.

Retenons donc de tout cela que l'estimateur sans biais de $g(\theta)$ et de variance minimale est une unique fonction d'une statistique exhaustive et complète. On aimerait maintenant connaître un peu mieux la covariance d'un tel estimateur.

Théorème (Lehmann-Scheffé) Soit T un estimateur sans biais de $g(\theta)$ et soit S une statistique exhaustive et complète. Soit $(\Omega, \mathcal{A}, \mathbb{P}_{\theta}, \theta \in \Theta)$ dominé et régulier, et soit T un estimateur sans biais de $g(\theta)$, tel que $E_{\theta}(T) = g(\theta)$. Si on suppose que l'information de Fisher est une matrice définie positive, alors, en notant $\frac{\partial g}{\partial \theta}(\theta)$ la matrice jacobienne de g , pour tout $\theta \in \Theta$:

$$\text{cov}(T) \geq \frac{\partial g}{\partial \theta}(\theta) \cdot (I_n(\theta))^{-1} \cdot \frac{\partial g}{\partial \theta}(\theta) \quad (\text{au sens des matrices symétriques}).$$

En particulier, si T est un estimateur sans biais de θ , alors:

$$\text{cov}(T) \geq (I_n(\theta))^{-1} \quad (\text{au sens des matrices symétriques}).$$

Démonstration Soit $Z_{\theta}(x) = \frac{\partial \log L_{\theta}(x)}{\partial \theta}$ où $x \in (\Omega')^n$ suit \mathbb{P}_{θ} . On sait que comme le modèle est régulier $E_{\theta}(Z_{\theta}) = 0$ pour tout $\theta \in \Theta$ et donc:

$$\text{cov}(Z_{\theta}) = I(\theta) \quad \text{pour tout } \theta \in \Theta.$$

De plus, T est un estimateur sans biais de $g(\theta)$ donc pour tout $\theta \in \Theta$

$$\begin{aligned} E_{\theta}(T) &= g(\theta) \implies \int_{(\Omega')^n} T(x) \cdot \frac{\partial L_{\theta}}{\partial \theta}(x) d\mu(x) = \frac{\partial g}{\partial \theta}(\theta) \quad (\text{en dérivant}) \\ &\implies \int_{(\Omega')^n} T(x) \cdot \frac{\partial L_{\theta}}{\partial \theta}(x) \cdot (L_{\theta}(x))^{-1} d\mathbb{P}_{\theta}(x) = \frac{\partial g}{\partial \theta}(\theta) \\ &\implies E_{\theta}(T \cdot Z_{\theta}) = \frac{\partial g}{\partial \theta}(\theta). \end{aligned}$$

Ainsi, d'après ce que l'on a vu,

$$\begin{aligned} \text{cov}\left(T - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta\right) &= \text{cov}_\theta(T) - 2 \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) + \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) \\ &= \text{cov}_\theta(T) - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta). \end{aligned}$$

En conséquence, comme $\text{cov}\left(T - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta\right)$ est une matrice positive définie, l'addition de Cramer-Rao est prouvée. ■

Corollaire Deux cas particuliers intéressants :

- Si le modèle est de la forme $(\Omega, \mathcal{A}_n, (f_\theta \cdot d\mu)^{\otimes n}, \theta \in \Theta)$, alors $s_{I_n}(\theta) = n \cdot I_1(\theta)$, où $I_1(\theta)$ est la matrice d'information de Fisher d'une seule variable aléatoire f_θ et l'égalité de Cramer-Rao devient donc :

$$\text{cov}(T) \geq \frac{1}{n} \cdot \frac{\partial g}{\partial \theta}(\theta) \cdot (I_1(\theta))^{-1} \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) \quad (\text{au sens des matrices symétriques}).$$

On voit donc qu'en augmentant le nombre de variables indépendantes, si la vraisemblance est régulière, alors la vitesse de convergence de tout estimateur sans biais est au mieux en $1/n$.

- Si le modèle n'est pas régulier mais qu'il est sous la probabilité P_θ , la matrice d'information de Fisher existe et est inversible. La propriété (4) est vérifiée, alors l'égalité de Cramer-Rao est vérifiée. Cela exclut cependant les modèles où les supports de f_θ dépendent de θ , comme par exemple le simple échantillon i.i.d. de $\text{Exp}([0, \theta])$, avec $\theta > 0$.

Définition Si un estimateur sans biais atteint (respectivement asymptotiquement) la borne de Cramer-Rao (qui ne dépend pas de l'estimateur), on dit qu'il est (asymptotiquement) efficace.

Remarque Un estimateur sans biais de variance minimale n'existe pas toujours. On ne peut pas atteindre la borne de Cramer-Rao, donc n'est pas efficace. De même, il peut exister des estimateurs sans biais n'atteignant pas la borne de Cramer-Rao.

Nous allons voir que les modèles exponentiels jouent un rôle central pour l'estimation sans biais efficace. Nous verrons que sous certaines conditions ils sont les seuls pour lesquels on a une estimation sans biais efficace.

Théorème Soit un modèle statistique paramétrique $(\Omega^n, \mathcal{A}_n, P_\theta, \theta \in \Theta)$, avec $\Theta \subset \mathbb{R}^p$, dont on suppose que la matrice de covariance $\frac{\partial g}{\partial \theta}(\theta)$ soit de rang plein pour tout $\theta \in \Theta$. Alors $T = (T_1, \dots, T_p)$ est un estimateur sans biais de $g(\theta)$ atteignant la borne de Cramer-Rao si et seulement si le modèle est exponentiel. Plus précisément, il existe des fonctions $\alpha_j : \Theta \rightarrow \mathbb{R}$, $\beta : \Theta \rightarrow \mathbb{R}$ et $\alpha_j : \Theta \rightarrow \mathbb{R}$ ($1 \leq j \leq p$), telles que pour tout $\theta \in \Theta$, $g(\theta) = \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \frac{\partial \beta}{\partial \theta}(\theta)$ et

$$L_\theta(x_1, \dots, x_n) = \exp\left[\beta(\theta) + \sum_{j=1}^p T_j(x_1, \dots, x_n) \cdot \alpha_j(\theta)\right].$$

Démonstration On suppose donc le modèle exponentiel. Si on écrit par rapport à θ un tel modèle, on obtient que pour μ -presque tout x

$$\frac{\partial}{\partial \theta}(\log L_\theta(x)) = \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot T_j + \frac{\partial \beta}{\partial \theta}(\theta), \quad \text{pour tout } \theta \in \Theta. \quad (5)$$

En conséquence, comme $I(\theta) = E\left[\frac{\partial}{\partial \theta}(\log L_\theta(.)) \cdot {}^t \frac{\partial}{\partial \theta}(\log L_\theta(.))\right]$, on en déduit que :

$$I(\theta) = \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \text{cov}(T) \cdot {}^t \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \quad \Rightarrow \quad \text{cov}(T) = \frac{\partial \alpha_j}{\partial \theta_i}(\theta)^{-1} \cdot I(\theta) \cdot {}^t \frac{\partial \alpha_j}{\partial \theta_i}(\theta)^{-1}$$

Par ailleurs, comme T est un estimateur sans biais de $g(\theta)$ la preuve de l'égalité de Cramér-Rao,

$$E_{\theta} \left[T(\cdot) \cdot \frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot)) \right] = \frac{\partial g}{\partial \theta}(\theta)$$

et en utilisant (5) que l'on multiplie par $\frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot))$, on obtient :

$$E_{\theta} \left[\frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot)) \cdot \frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot)) \right] = E_{\theta} \left[\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot T \cdot \frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot)) + \frac{\partial \beta}{\partial \theta}(\theta) \cdot \frac{\partial}{\partial \theta} (\log l_{\theta}(\cdot)) \right],$$

et donc $I(\theta) = \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \frac{\partial g}{\partial \theta}(\theta)$. A l'aide de cette égalité en reprenant le calcul précédent, on en arrive à ce que :

$$\text{cov}(T) = \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot \frac{\partial g}{\partial \theta}(\theta),$$

donc T atteint bien la borne de Cramér-Rao. De plus, grâce à (5),

$$E_{\theta} \left[\frac{\partial}{\partial \theta} (\log l_{\theta}(x)) \right] = E_{\theta} \left[\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot T + \frac{\partial \beta}{\partial \theta}(\theta) \right]$$

$$\text{soit} \quad 0 = \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot g(\theta) + \frac{\partial \beta}{\partial \theta}(\theta)$$

$$\text{et donc} \quad g(\theta) = - \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \frac{\partial \beta}{\partial \theta}(\theta).$$

⇒ D'après la preuve de l'égalité de Cramér-Rao, si T est un estimateur sans biais de $g(\theta)$ atteignant la borne de Cramér-Rao, alors

$$\text{cov}(T - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_{\theta}) = 0.$$

Ainsi, pour tout $\theta \in \Theta$, il existe un ensemble $N(\theta)^n$ tel que $P(N_{\theta}) = 1$ et tel que pour tout $\omega \in N_{\theta}$,

$T(x) - g(\theta) = \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_{\theta}(x)$. Par le même principe que celui de la preuve de la nullité de la matrice de Fisher pour une statistique libre, on peut trouver un ensemble N dépendant pas de θ , tel que cette propriété se vérifie également, avec $P(N) = 1$, ce qui implique que :

$$I(\theta) \cdot \frac{\partial g}{\partial \theta}(\theta)^{-1} \cdot T(x) - g(\theta) = \frac{\partial}{\partial \theta} (\log l_{\theta}(x)), \quad \text{pour tout } \theta \in \Theta.$$

Alors en intégrant par rapport à θ en notant :

- $\alpha(\theta)$ le vecteur colonne "grand" $I(\theta) \cdot \frac{\partial g}{\partial \theta}(\theta)^{-1}$
- $\beta(\theta)$ la fonction "intégrale" $\frac{\partial g}{\partial \theta}(\theta)^{-1} \cdot g(\theta)$
- $b(x)$ une fonction ne dépendant pas de θ

on a $\log l(x) = \alpha(\theta) \cdot T(x) + \beta(\theta) + b(x)$, d'où l'écriture de la vraisemblance sous forme canonique d'expérimentiel, et on retrouve l'expression de $g(\theta)$ par conséquent plus haut. ■

Corollaire : l'inverse si l'on dispose d'un modèle exponentiel régulier (3), alors il n'existe qu'une seule fonction (à un étirement près) afin d'appréhender le pourcentage d'efficacité d'estimation

$$g(\theta) = \frac{1}{n} \cdot \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \frac{\partial \beta}{\partial \theta}(\theta) \quad (\text{noter que cette fonction se simplifie de } n; \text{ dans le cas de v.a.i.i.d.}$$

ce n'est pas le cas). L'estimateur est alors $\frac{1}{n} \cdot T = \frac{1}{n} \cdot (X_1, \dots, X_n, \dots, X_1, \dots, X_n)$ et sa matrice de covariance minimale est donc par sa borne de Cramér-Rao, soit :

$$\text{cov}(T) = \frac{1}{n} \cdot \frac{\partial g}{\partial \theta}(\theta) \cdot \frac{\partial \alpha_j}{\partial \theta_i}(\theta)^{-1}.$$

3.5 Estimateur du maximum de vraisemblance

Nous allons voir une méthode permettant d'obtenir aisément et dans la plupart des cas un estimateur possédant de très bonnes propriétés. Par la suite on se place une nouvelle fois dans le cadre d'un modèle statistique paramétrique $((\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta))$, avec $\Theta \subset \mathbb{R}^p$, comme d'habitude.

Définissons pour $(x_1, \dots, x_n) \in (\Omega)^n$, soit $\theta \in \Theta \rightarrow L_\theta(x_1, \dots, x_n)$ la vraisemblance du échantillon. On appelle estimateur du maximum de vraisemblance un estimateur $\hat{\theta}$ pour (x_1, \dots, x_n) un n-échantillon quelconque du échantillon.

$$L_{\theta_n}(X_1, \dots, X_n) = \sup_{\theta \in \Theta} L_\theta(X_1, \dots, X_n).$$

Remarquons qu'il n'y a pas de garantie de l'unicité d'un tel estimateur. Un exemple pour l'obtenir (mais pas toujours) est de rechercher un extremum local sur Θ , ce qui peut se faire annuellement dérivées partielles de $\log L_\theta$. De même, il est clair que l'estimateur du maximum de vraisemblance pour θ est également obtenu en maximisant le logarithme de la vraisemblance, appelée vraisemblance. Enfin, si l'on désire estimer $g(\theta)$ avec g une fonction bijective, alors $g(\hat{\theta})$ sera l'estimateur du maximum de vraisemblance de $g(\theta)$.

Propriété. S'il existe une statistique exhaustive T pour θ , alors $\hat{\theta}$ est une fonction mesurable de T pour tout $\theta \in \Theta$.

Démonstration: Si T est exhaustive, d'après le théorème de factorisation, la vraisemblance se factorise en $L_\theta(x_1, \dots, x_n) = h_\theta(T(x_1, \dots, x_n)) \cdot k_\theta(x_1, \dots, x_n)$ pour tout $\theta \in \Theta$ et \mathbb{P}_θ -presque tout $(x_1, \dots, x_n) \in (\Omega)^n$, ce qui revient à dire que presque tout $(x_1, \dots, x_n) \in (\Omega)^n$ par la même démonstration que celle de la nullité de l'information de Fisher d'une statistique exhaustive. Ainsi, pour prendre l'argument maximal de L_θ revient à prendre l'argument maximal de $h_\theta(T(x_1, \dots, x_n))$, et $\hat{\theta}$ sera donc une fonction de T . ■

Propriété. On suppose que le modèle est régulier. Si on suppose qu'il existe un estimateur sans biais efficace de θ alors c'est l'estimateur du maximum de vraisemblance de θ .

Démonstration: D'après ce qui précède, si le modèle est régulier et que T est un estimateur sans biais efficace de θ , alors le moment des T est égal à θ (5) à l'encore lieu, soit pour tout $\theta \in \Theta$

$$\frac{\partial}{\partial \theta} (\log L_\theta(x)) = \sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot T + \frac{\partial \beta}{\partial \theta}(\theta) \Rightarrow \sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \mathbb{E}_\theta(T) + \frac{\partial \beta}{\partial \theta}(\theta) = 0.$$

Comme T est un estimateur sans biais de θ , on a donc $\sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\theta) \cdot \theta + \frac{\partial \beta}{\partial \theta}(\theta) = 0$, pour tout $\theta \in \Theta$, ce qui s'applique également à $\hat{\theta}$ et donc:

$$\sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\hat{\theta}) \cdot \hat{\theta} + \frac{\partial \beta}{\partial \theta}(\hat{\theta}) = 0.$$

Mais d'après la définition, le moment des $\hat{\theta}$ est égal à $\hat{\theta}$ et $\hat{\theta}$ minimise la log-vraisemblance et a donc une variance nulle, ce qui implique que:

$$\sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\hat{\theta}) \cdot T + \frac{\partial \beta}{\partial \theta}(\hat{\theta}) = 0.$$

En conséquence, on obtient:

$$\sum_{1 \leq i, j \leq p} \frac{\partial \alpha_j}{\partial \theta_i}(\hat{\theta}) \cdot T - \hat{\theta} = 0 \Rightarrow T = \hat{\theta},$$

car la matrice des dérivées secondes est supposée de rang d. Enfin, l'unicité de l'écriture du moment est évidente. ■

Nous allons nous intéresser maintenant au comportement asymptotique de l'estimateur du maximum de vraisemblance (loisqu'il existe), donc quand la taille de l'échantillon tend vers l'infini. Il est clair que pour chaque n l'expression de l'estimateur est explicite, toute fois, pour le modèle statistique cherché par ailleurs, cela, on se placera dans un "gros" modèle $((\Omega, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta))$, où $\Theta \subset \mathbb{R}^p$ (la dimension du paramètre n'est pas forcément finie) dans lequel on suppose que $\theta \mapsto L_\theta(x_1, \dots, x_n)$ est une suite de fonctions continues et différentiables. Par ailleurs, on suppose aussi que tout échantillon de ce modèle est constitué de v. a. i. i. d., et que $\mathbb{P}_\theta = \mu^{\otimes n}$, le modèle étant donc par la mesure μ et μ étant la densité cherchée par rapport à μ .

Théorème (Convergence de l'estimateur du maximum de vraisemblance) On suppose le modèle $((\Omega)^{\mathbb{N}}, \mathcal{A}_{\mathbb{N}}, (f_{\theta} \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\theta \in \mathbb{R}^d$ domine par une mesure régulière. On suppose en plus que le modèle est identifiable (au sens de $f_{\theta_1} \neq f_{\theta_2}$ presque partout si $\theta_1 \neq \theta_2$). Alors si la suite $(X_n)_{n \in \mathbb{N}}$ est issue du modèle avec pour paramètre $\theta \in \Theta$,

$$\theta_n \xrightarrow[n \rightarrow +\infty]{p.s.} \theta \quad \text{pour la mesure } (f_{\theta} d\mu)^{\otimes \mathbb{N}}.$$

Démonstration: En premier lieu, on fixe et on se propose de montrer que pour tout $\theta \in \Theta$

$$\log \ell(X_1, \dots, X_n) - \log \ell_{\theta}(X_1, \dots, X_n) = \sum_{i=1}^n \log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)}.$$

Par ailleurs, pour tout $\theta \in \Theta$, les X_i ont tous la même loi et pour tout $\theta \in \Theta$,

$$\begin{aligned} \mathbb{E}_{\theta_0} \log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} &\leq \log \mathbb{E}_{\theta_0} \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \quad (\text{Inégalité de Jensen pour la fonction log}) \\ &\leq \log (\mathbb{E}[f_{\theta}(X_i)]) \\ &\leq 0. \end{aligned}$$

En fait, du fait que la fonction log est strictement concave, la borne atteinte que si $\theta = \theta_0$. Ainsi, avec la contrainte d'un paramètre identifiable d'où $\theta \neq \theta_0$ on a :

$$\mathbb{E}_{\theta_0} \log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} < 0.$$

On peut appliquer la loi forte des grands nombres pour les variables $\log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)}$ (qui sont bien i.i.d. et d'après le modèle est régulier), et ainsi :

$$\begin{aligned} \frac{1}{n} (\log \ell(X_1, \dots, X_n) - \log \ell_{\theta_0}(X_1, \dots, X_n)) &= \frac{1}{n} \sum_{i=1}^n \log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \\ &\xrightarrow[n \rightarrow +\infty]{p.s.} \mathbb{E}_{\theta_0} \log \frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} < 0, \end{aligned}$$

la convergence presque sûre ayant lieu pour la mesure $(f_{\theta_0} d\mu)^{\otimes \mathbb{N}}$ (il nous faut donc maintenant pour tout $\varepsilon > 0$ une famille nombreable $\mathcal{I}_{\varepsilon}$ dense sur la boule de centre θ et de rayon ε . Du fait du caractère dénombrable de cette famille, pour tout $\varepsilon > 0$, il existe pour tout $\theta \in \Theta$:

$$\log \ell_{\theta_1}(X_1, \dots, X_n) < \log \ell_{\theta}(X_1, \dots, X_n) \quad \text{p.s. pour la mesure } (f_{\theta} d\mu)^{\otimes \mathbb{N}}.$$

Comme le modèle est régulier, pour tout $\theta \in \Theta$, la log-vraisemblance de X, X est continue sur Θ . De plus pour tout n elle atteint son unique maximum en θ_n et donc, pour une θ se rapprochant de la boule de centre θ de rayon ε (toujours p.s. pour la mesure $(f_{\theta} d\mu)^{\otimes \mathbb{N}}$). Le raisonnement est valable pour tout $\varepsilon > 0$, et on s'en suit. ■

Théorème (Normalité asymptotique de l'estimateur du maximum de vraisemblance) On suppose le modèle paramétrique $((\Omega)^{\mathbb{N}}, \mathcal{A}_{\mathbb{N}}, (f_{\theta} \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$ domine par une mesure régulière. On suppose en plus que le modèle est identifiable et que la fonction $\theta \mapsto L_{\theta}$ est de classe $\mathcal{C}^2(\Theta)$. Alors si la suite $(X_n)_{n \in \mathbb{N}}$ est issue du modèle avec pour paramètre $\theta \in \Theta$:

$$\sqrt{n}(\theta_n - \theta) \xrightarrow[n \rightarrow +\infty]{L} N_d(0, I^{-1}(\theta)),$$

où $I(\theta)$ est la matrice de Fisher de taille p (supposable) pour la variable X

Démonstration: Comme le modèle est régulier, on peut différencier la vraisemblance et pour tout $\theta \in \Theta$ on a :

$$M_{\theta}(X_1, \dots, X_n) = \frac{1}{n} \frac{\partial}{\partial \theta} \log \ell_{\theta}(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log \ell(X_i).$$

Un développement limité d'ordre 1 de M_{θ} autour de θ_0 est possible (toujours en raison d'une dérivée existante) et donc possible tout à fait :

$$M_{\theta}(X_1, \dots, X_n) = M_{\theta_0}(X_1, \dots, X_n) + (\theta - \theta_0) \cdot \frac{\partial}{\partial \theta} M_{\theta_0^*}(X_1, \dots, X_n),$$

avec dans le segment $[\theta_0, \theta]$, on remarque que $\frac{\partial}{\partial \theta} M_{\theta_0^*}(X_1, \dots, X_n)$ est une matrice carrée (taillée). Ainsi en remplaçant θ_0 par θ , on obtient pour chaque n l'existence d'un θ_n tel que :

$$M_{\theta_n}(X_1, \dots, X_n) = M_{\theta_0}(X_1, \dots, X_n) + (\theta_n - \theta_0) \cdot \frac{\partial}{\partial \theta} M_{\theta_n^*}(X_1, \dots, X_n). \quad (6)$$

Pour un modèle régulier, on a vu que $\frac{\partial^2}{\partial \theta^2} \log f_{\theta_0}(X_i) = -I_1(\theta_0)$, matrice de Fisher pour n'importe quelle variable. Ainsi, $\frac{\partial}{\partial \theta} M_{\theta}(\cdot)$ étant une moyenne empirique, on a par la loi forte des grands nombres :

$$\frac{\partial}{\partial \theta} M_{\theta_0}(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f_{\theta_0}(X_i) \xrightarrow[n \rightarrow +\infty]{p.s.} -I_1(\theta_0) \text{ pour la mesure } (\mathbb{P}_{\mu})^{\otimes \mathbb{N}}.$$

Maintenant, en utilisant le fait que les dérivées de classe $\mathcal{C}^2(\Theta)$ et en utilisant la convergence presque sûre de $\frac{\partial}{\partial \theta} M_{\theta_n^*}(X_1, \dots, X_n)$ à $-I_1(\theta_0)$, on a :

$$\frac{\partial}{\partial \theta} M_{\theta_n^*}(X_1, \dots, X_n) \xrightarrow[n \rightarrow +\infty]{p.s.} -I_1(\theta_0) \text{ pour la mesure } (\mathbb{P}_{\mu})^{\otimes \mathbb{N}}.$$

Finalement, comme c'est le maximum d'une fonction de classe \mathcal{C}^1 , le θ_n est un estimateur annulé $M_{\theta_n}(X_1, \dots, X_n)$, et donc l'équation (6) devient :

$$M_{\theta_0}(X_1, \dots, X_n) \cdot I_1^{-1}(\theta_0) = (\theta_n - \theta_0).$$

Enfin, comme $M_{\theta_0}(X_1, \dots, X_n)$ est une moyenne empirique, avec une attente vérifiant un théorème de la limite centrale :

$$\sqrt{n} (M_{\theta_0}(X_1, \dots, X_n) - \mathbb{E}_{\theta_0} \frac{\partial}{\partial \theta} \log f_{\theta_0}(X_i)) \xrightarrow[n \rightarrow +\infty]{L} N_d(0, I_1(\theta_0)),$$

d'où par la première définition de l'information de Fisher, on a $\frac{\partial}{\partial \theta} \log f_{\theta_0}(X_i) = 0$ (voir les propriétés précédentes), on obtient la notation asymptotique de θ_n : $\theta_n \sim \theta_0 + \sqrt{n} (M_{\theta_0}(X_1, \dots, X_n) - \mathbb{E}_{\theta_0} \frac{\partial}{\partial \theta} \log f_{\theta_0}(X_i)) I_1^{-1}(\theta_0)$. ■

Remarque : Sous ces hypothèses, l'estimateur du maximum de vraisemblance est asymptotiquement biaisé et efficace. Cependant, il peut avoir un biais et ne pas être un estimateur efficace.

3.6 Régions de confiance

En pratique, estimer un paramètre plus souvent ne suffit pas. On aiment le plus précisément quelle marge de confiance a sur la connaissance de ce paramètre :

Définition : On se place dans le cadre d'un modèle paramétrique $((\Omega, \mathcal{A}_n, \mathbb{P}_{\theta}, \theta \in \Theta))$, où $\Omega \subset \mathbb{R}^p$. Soit $\alpha \in]0, 1[$ un nombre fixé a priori. On appelle région de confiance du paramètre de niveau $1-\alpha$ un sous-ensemble $R_{1-\alpha}$ inclus dans \mathbb{R}^p et défini sur $((\Omega, \mathcal{A}_n))$, tel que pour tout $\theta \in \Theta$, $\{(X_1, \dots, X_n) \in (\Omega)^n, \theta \in R_{1-\alpha}(X_1, \dots, X_n)\} \in \mathcal{A}_n$ et :

$$\inf_{\theta \in \Theta} \{\mathbb{P}_{\theta}(\theta \in R_{1-\alpha})\} \geq 1 - \alpha. \quad (7)$$

Si un échantillon observé $(X_1(\omega), \dots, X_n(\omega))$ est connu, $R_{1-\alpha}(X_1(\omega), \dots, X_n(\omega))$ est appelée région de confiance observée. Dans le cas où le paramètre est un réel ($p = 1$), on pourra obtenir un intervalle de confiance.

Comment évaluer une région de confiance ? En premier lieu, il est clair que $\theta_0 \in \Theta$, si $\theta_0 \in \Theta$ (en général, on choisit α proche de 0, et en particulier $\alpha = 0.05$). On est alors dans une situation possible pour la construction de région de confiance est la suivante : nous aurons besoin d'utiliser un estimateur T convergent de θ , mais s'agit-il de θ ce qui rend difficile (quelques exceptions) son utilisation directe. On préférera donc utiliser ce que l'on appelle une fonction pivotale qui est θ une fonction mesurable d'un estimateur et de cette statistique libre. On essaiera alors d'écrire la probabilité (7) sous la forme

$$\inf_{\theta \in \Theta} P_{\theta}(\pi(T, \theta) \in C_{\alpha}) \geq 1 - \alpha,$$

où C_{α} est une région de confiance minimale. On s'occupera ensuite de construire la région de confiance en fonction des quantiles (souvent à $\alpha/2$ et $1-\alpha/2$) de la loi de la fonction pivotale.

Exemple : Si le modèle est régulier, sous les conditions de normalité asymptotique du maximum de vraisemblance, on peut même montrer (voir le théorème de Slutski) que

$$\pi(\theta_n, \theta_0) = \sqrt{n} \cdot (I_1(\theta_n))^{1/2} \cdot \theta_n - \theta_0 \xrightarrow[n \rightarrow +\infty]{L} N_d(0, b),$$

où b est la matrice identité et $I_1(\theta) = \frac{1}{n} \cdot (I_1(\theta))^{1/2} = I_1(\theta)$ pour tout $\theta \in \Theta$. Ainsi, si n est grand, on pourra assimiler la loi de $\pi(\theta_n, \theta_0)$ avec la loi normale centrée multidimensionnelle. Or si $Z \sim N_p(0, b)$, avec $q_{1-\alpha/2} > 0$ le quantile d'une normale centrée réduite de niveau $1-\alpha/2$, tel que $P(Z \in [-q_{1-\alpha/2}, q_{1-\alpha/2}]^p) \geq 1 - \alpha$. Aussi le pivot $\sqrt{n} \cdot (I_1(\theta_n))^{-1/2} \cdot [-q_{1-\alpha/2}, q_{1-\alpha/2}]^d$ recouvrera l'ensemble de la région de confiance cherchée.

3.7 M-estimateur

Il s'agit ici d'évaluer les estimateurs du maximum de vraisemblance. On donne ainsi des critères généraux pour la consistance et l'asymptotique des estimateurs.

3.7.1 Introduction

Supposons que nous voulons estimer θ relié à la loi de probabilité (X_n) . La méthode pour trouver un tel estimateur est de minimiser une fonctionnelle

$$\theta \rightarrow M_n(\theta) = \frac{1}{n} \sum_{i=1}^n m_{\theta}(X_i)$$

avec $m_{\theta} : \mathbb{R}^p \rightarrow \mathbb{R}$ sont des fonctions continues.

Exemples : estimateur du maximum de vraisemblance. Admettons que (X_1, \dots, X_n) sont i.i.d. leur loi à θ pour lequel l'estimateur du maximum de vraisemblance est :

$$\theta \rightarrow M_n(\theta) = \frac{1}{n} \sum_{i=1}^n -\log f_{\theta}(X_i)$$

Exemple : estimateur des moindres carrés. Admettons que (X_1, \dots, X_n, Y_n) sont i.i.d., $X \in \mathbb{R}^p, Y_i \in \mathbb{R}$ et vérifient l'équation de régression linéaire :

$$Y = \theta^T X + \varepsilon$$

où ε est une variable aléatoire indépendante de X et de θ et ε est centrée. L'estimateur des moindres carrés $\hat{\theta}_n$ de θ est alors celui qui minimise :

$$\theta \rightarrow M_n(\theta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \theta^T X_i)^2$$

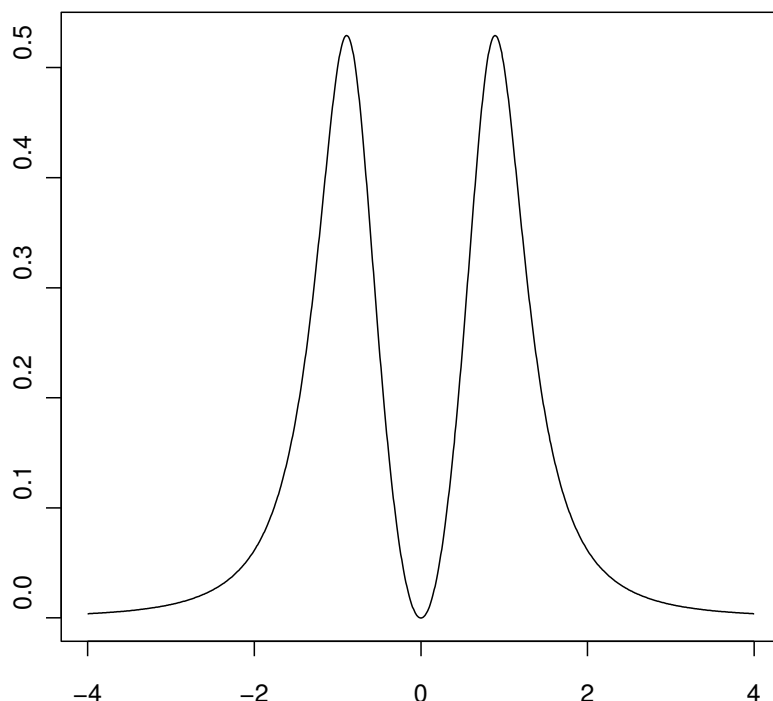


Figure 1 - Minimum multiple

3.7.2 Consistance

Il est important que l'estimateur converge vers la vraie valeur θ_0 presque sûrement ou bien en probabilité. Il faut donc que le nombre d'observations n converge vers l'infini. Si c'est le cas, l'estimateur est dit asymptotiquement consistant. Par exemple la moyenne \bar{X}_n est asymptotiquement consistante pour la moyenne de la population $\theta = E(X)$ si $E(X)$ existe. On veut donc prouver que:

$$\hat{\theta}_n \xrightarrow{p.s.} \theta$$

On suppose que le M-estimateur minimise la fonction $M_n(\theta)$. Clairement, le comportement asymptotique de $\hat{\theta}_n$ dépend du comportement asymptotique de la fonction M_n . On suppose qu'il existe une fonction réelle $M(\theta)$ telle que:

$$\forall \theta : M_n(\theta) \xrightarrow{p.s.} M(\theta)$$

Il semblerait raisonnable que le minimum de $M_n(\theta)$ converge sous des conditions raisonnables vers le minimiseur de $M(\theta)$. Il faut quand même faire attention à ce que le problème de minimisation de $M(\theta)$ soit bien posé. Éviter ce genre de situation. La figure 1 donne un exemple d'une fonction qui a bien un minimum égal à 0, mais la droite $y = 0$ est aussi une asymptote.

Théorème Soit Θ l'espace des paramètres possibles et M des fonctions telles que pour tout $\varepsilon > 0$:

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{p.s.} 0, \\ \inf_{\theta : |\theta - \theta_0| \geq \varepsilon} M(\theta) > M(\theta_0)$$

alors pour toutes suites d'estimateurs $\hat{\theta}_n$ on a que

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} M_n(\theta)$$

On aura

$$\hat{\theta}_n \xrightarrow{p.s.} \theta_0$$

Démonstration On a pour tout voisinage V de θ_0 l'existence d'une constante $\eta(V)$ telle que :

$$\forall \theta \in \Theta \setminus V, M(\theta) > M(\theta_0) + \eta(V)$$

Donc, pour montrer la consistance forte, il suffit de montrer que $\lim_{n \rightarrow \infty} M(\hat{\theta}_n) = M(\theta_0)$ p.s.

$$\lim_{n \rightarrow \infty} \hat{\theta}_n \subset V \Leftrightarrow \lim_{n \rightarrow \infty} M(\hat{\theta}_n) - M(\theta_0) \leq \eta(V) \text{ p.s.}$$

Par définition $M(\hat{\theta}_n) \leq M_n(\hat{\theta}_n)$ et comme $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{p.s.} 0$ on aura

$$\lim_{n \rightarrow \infty} M_n(\hat{\theta}_n) \leq M(\theta_0) + \frac{\eta(V)}{2} \text{ p.s.}$$

De même on a

$$\lim_{n \rightarrow \infty} M(\hat{\theta}_n) - M_n(\hat{\theta}_n) \xrightarrow{p.s.} 0$$

et

$$\lim_{n \rightarrow \infty} M(\hat{\theta}_n) - \frac{\eta(V)}{2} \leq \lim_{n \rightarrow \infty} M_n(\hat{\theta}_n) \leq M(\theta_0) + \frac{\eta(V)}{2} \text{ p.s.}$$

finalement $\lim_{n \rightarrow \infty} M(\hat{\theta}_n) - M(\theta_0) \leq \eta(V)$ p.s. ce qui prouve la consistance forte du M-estimateur.

Conditions suffisantes pour la consistance forte. Pour un modèle régulier les hypothèses du théorème sont faciles à vérifier. Pour obtenir la condition :

$$\theta: \inf_{\theta - \theta_0 \geq \varepsilon} M(\theta) > M(\theta_0)$$

Il suffit que la fonction limite $M(\theta)$ soit une fonction strictement minimale en θ_0 et que pour tout $\theta \neq \theta_0$, $M(\theta) = M(\theta_0)$. Cela se traduit par le fait que $M(\theta)$ est l'espérance de l'opposé de la log-vraisemblance, c'est-à-dire la distance de Kullback à une constante pr

$$M(\theta) = -f_{\theta_0} \log \frac{f_{\theta_0}}{f_{\theta}}$$

Pour obtenir l'hypothèse de la loi uniforme des grands nombres :

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{p.s.} 0$$

une condition suffisante est que l'ensemble des paramètres possibles Θ soit compact, la fonction $m_{\theta}(x)$ soit continue pour tout x et qu'il existe une fonction $m(x)$ qui do

Preuve Pour une boule ouverte B de Θ on a $m_B = \sup_{\theta \in B} m_{\theta}$ et $m_B = \inf_{\theta \in B} m_{\theta}$. Par le lemme de convergence de $E[m_B] - m_B \rightarrow 0$ lorsque le diamètre de la boule tend vers 0. Pour $\varepsilon > 0$, soit B^1, \dots, B^k un recouvrement fini de Θ tel que $E[m_{B^i}(X)] < \varepsilon$. Pour tout $\theta \in B^i$, on a :

$$\begin{aligned} M_n(\theta) - M(\theta) &\leq \frac{1}{n} m^{B^i}(X_i) - E[m^{B^i}(X)] \leq \frac{1}{n} m^{B^i}(X_i) - E[m^{B^i}(X)] + \varepsilon \\ M_n(\theta) - M(\theta) &\geq \frac{1}{n} m_{B^i}(X_i) - E[m^{B^i}(X)] \geq \frac{1}{n} m_{B^i}(X_i) - E[m_{B^i}(X)] - \varepsilon \end{aligned}$$

Ainsi,

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \leq \sup_{i \in \{1, \dots, k\}} \max \left(\frac{1}{n} m^{B^i}(X_i) - E[m^{B^i}(X)], \frac{1}{n} m_{B^i}(X_i) - E[m_{B^i}(X)] \right) + \varepsilon$$

Ainsi, presque sûrement, $\lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| < \varepsilon$ ■

3.7.3 Normalité Asymptotique

On donne ici des conditions nécessaires et suffisantes pour la normalité asymptotique du M-estimateur. Ces conditions peuvent être affaiblies mais celles données conviennent pour les exemples considérés.

Notations

- Pour une suite de variables X_n on écrit

$$X_n = o_P(R_n) \Leftrightarrow X_n = Y_n R_n \text{ et } Y_n \xrightarrow{P} 0$$

- Pour une suite de variables X_n on écrit

$$X_n = O_P(R_n) \Leftrightarrow X_n = Y_n R_n \text{ et } Y_n \text{ est bornée en probabilité (il existe } M \text{ tel que } P(|Y_n| > M) \rightarrow 0).$$

Hypothèses On suppose que les observations (X_n) sont i.i.d. et que les hypothèses suivantes sont vérifiées:

H1 Le M-estimateur est fortement consistant

H2 Il existe un voisinage V du vrai paramètre tel que pour tout $\theta \in V$, la dérivée $\frac{\partial}{\partial \theta} m_\theta(X)$ est dominée par une fonction intégrable.

H3 Le carré de la dérivée $\frac{\partial m_\theta(X)}{\partial \theta}(\theta_0)^2$ est intégrable.

H4 La matrice des dérivées secondes $\frac{\partial^2 m_\theta(X)}{\partial \theta^2}(\theta_0)$ est intégrable et inversible.

On a alors la conséquence suivante:

Théorème Sous les hypothèses H1, ..., H4:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = -E\left[\frac{\partial^2 m_\theta}{\partial \theta^2}(\theta_0)\right]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial m_\theta}{\partial \theta}(\theta_0) + o_P(1)$$

En particulier la suite $\sqrt{n}(\hat{\theta}_n - \theta_0)$ est asymptotiquement normale de moyenne 0 et de matrice de covariance $E\left[\frac{\partial^2 m_\theta}{\partial \theta^2}(\theta_0)\right]^{-1} E\left[\frac{\partial m_\theta}{\partial \theta}(\theta_0) \frac{\partial m_\theta}{\partial \theta}(\theta_0)^T\right] E\left[\frac{\partial^2 m_\theta}{\partial \theta^2}(\theta_0)\right]^{-1}$.

Démonstration: Par un développement de Taylor il existe un vecteur $\tilde{\theta}$ sur le segment $[\theta_0, \hat{\theta}_n]$ tel que:

$$0 = \frac{\partial M_n(\hat{\theta}_n)}{\partial \theta} = \frac{\partial M_n(\theta_0)}{\partial \theta} + \frac{\partial^2 M_n(\theta_0)}{\partial \theta^2}(\hat{\theta}_n - \theta_0) + \frac{1}{2}(\hat{\theta}_n - \theta_0)^T \frac{\partial^3 M_n(\theta_0)}{\partial \theta^3}(\hat{\theta}_n - \theta_0)$$

Le premier terme à droite est la moyenne du vecteur $\frac{\partial M_n(\theta_0)}{\partial \theta}$ qui a pour espérance:

$$E\left[\frac{\partial M_n(\theta_0)}{\partial \theta}\right] = 0$$

Par le théorème de la limite centrale la suite $\frac{1}{\sqrt{n}} \frac{\partial M_n(\theta_0)}{\partial \theta}$ converge en loi vers une gaussienne de moyenne 0 et de matrice de variance-covariance $E\left[\frac{\partial m_\theta(\theta_0)}{\partial \theta} \frac{\partial m_\theta(\theta_0)}{\partial \theta}^T\right]$. Par la loi des grands nombres $\frac{1}{n} \frac{\partial^2 M_n(\theta_0)}{\partial \theta^2}$ converge

presque-sûrement vers une matrice $E\left[\frac{\partial^2 m_\theta(\theta_0)}{\partial \theta^2}\right] = J$. Si k est la dimension du vecteur θ , la dérivée

$\frac{\partial^3 M_n(\theta_0)}{\partial \theta^3}$ est un vecteur de k matrices k par hypothèse et existe un voisinage V de θ_0 tel que $\frac{\partial^3 M_n(\theta_0)}{\partial \theta^3}$

est dominé par une fonction intégrable. Comme $\hat{\theta}_n$ est consistant, $\lim_{n \rightarrow \infty} \hat{\theta}_n \in V$ presque-sûrement et pour $\hat{\theta}_n \in V$ on aura le développement suivant:

$$-\frac{\partial M_n(\theta_0)}{\partial \theta} = J + o_P(1) + \frac{1}{2}(\hat{\theta}_n - \theta_0)^T O_P(1) (\hat{\theta}_n - \theta_0) = (J + o_P(1))(\hat{\theta}_n - \theta_0)$$

car $(\hat{\theta}_n - \theta_0) O_P(1) = o_P(1)$ si $\hat{\theta}_n$ converge presque-sûrement vers θ_0 . On sait que la matrice $J + o_P(1)$ soit inversible tend vers 1. En multipliant membre à membre par $n(J + o_P(1))^{-1}$ on obtient le résultat annoncé. ■

4 Tests paramétriques

4.1 Principes d'un test

Un test permet, à partir d'un échantillon de données, de décider entre deux hypothèses, en mettant en avant une hypothèse privilégiée appelée H_0 et une hypothèse alternative appelée H_1 . On associe à un test un niveau α (avec $0 < \alpha < 1$) et une puissance β . La plupart du temps, est fixé a priori et β s'adapte. Plus précisément,

Définissons le cadre d'un problème statistique donné par $(\Omega, \mathcal{A}, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Omega \subset \mathbb{R}^p$ et soit θ la "vraie" valeur du paramètre. Un problème de test est un choix entre deux hypothèses

$$\begin{aligned} & \square H_0 : \theta \in \Theta_0 \quad : \text{hypothèse dite nulle} \\ & \square H_1 : \theta \in \Theta_1 \quad : \text{hypothèse dite alternative,} \end{aligned} \quad (8)$$

où $\Theta \subset \mathbb{R}^p, \Theta_1 \subset \mathbb{R}^d$ et $\Theta_0 \cap \Theta_1 = \emptyset$

Ceci peut se présenter deux types de problèmes de tests suivant les constitutions de Θ

Définissons une hypothèse (H_0 ou H_1) est dite simple si elle est associée à un singlet ($\theta \in \Theta_1$). Sinon, elle sera dite composite. Dans le cas ($\Theta \subset \mathbb{R}$), si H_0 est simple de la forme $\theta = \theta_0$ et si H_1 est composite de la forme $\theta > \theta_0$ ou $\theta < \theta_0$, on parlera de test unilatéral. Si H_1 est composite de la forme $\theta \neq \theta_0$ on parlera de test bilatéral.

Comment faire pour choisir entre les deux hypothèses H_0 et H_1 ? Il faudra partir de ce que l'on peut connaître du modèle c'est-à-dire généralement l'échantillon observé (X_1, \dots, X_n) . Pour cela on définit une statistique qui sera appelée T utile du test:

Dans le cadre du problème de test (8), soit T une statistique (donc une fonction mesurable d'un échantillon (X_1, \dots, X_n) issu du modèle) à valeurs dans \mathbb{R} qui sera appelée statistique du test. On définit la fonction $\phi = \mathbb{1}_{T \in W}$, où W est une partie de \mathbb{R} appelée région critique du test. La partie complémentaire dans \mathbb{R} sera appelée région d'acceptation du test. Si $\phi = 1$, on choisira H_1 on décidera plutôt H_0 .

Donc, à chaque hypothèse H_1 , on associe une partie de \mathbb{R} la statistique de test T et W . Ces parties ne sont pas disjointes. Pour pouvoir préciser la région W , dans un cadre continu (qui n'est pas le cas en pratique, voir plus bas), on peut commencer par associer une fonction puissance à la statistique de test T puis les règles de prise de décision de la façon suivante:

Définissons pour la statistique de test T , on associe:

- une fonction puissance, qui est la probabilité $\mathbb{P}_\theta(T \in W) : \theta \in \Theta_1 \rightarrow \mathbb{P}_\theta(T \in W)$.
- une erreur de première espèce : $\mathbb{P}_{H_0}(\text{Choisir } H_1) = \alpha = \sup_{\theta \in \Theta_0} \mathbb{P}_\theta(T \in W)$;
- une erreur de seconde espèce : $\mathbb{P}_{H_1}(\text{Choisir } H_0) = \beta = \sup_{\theta \in \Theta_1} \mathbb{P}_\theta(T \notin W)$.

La puissance du test est $1 - \beta$

Cependant, ce qui vient d'être écrit est technique. En pratique, on utilisera la méthode suivante:

Construction concrète d'un test: On suppose le problème de test (8). On pose également a priori α qui dépend du problème et $\alpha = 0.05$ et $1 - \alpha$ est appelé niveau du test. Par la suite, on procédera de la façon suivante:

1. L'expression quantitative des hypothèses H_0 et H_1 .
2. Le choix de la statistique T du test.
3. La construction d'une région critique W à l'aide de la statistique T .
4. La détermination explicite de W en fonction de α .
5. Le calcul (si possible) de la puissance du test $1 - \beta$.
6. Pour la réalisation de l'échantillon, rejeter ou accepter l'hypothèse H_0 .

Remarque : Cependant, pratiquement on ne peut pas s'inscrire dans ces deux types d'erreur. Le choix de l'hypothèse privilégiée est donc fondamentalement arbitraire. L'alternative d'un test n'est pas systématique. Par exemple, supposons que l'on ait pour $\theta \in \mathbb{R}$, $B(\mathbb{R}^n)$, $N(\theta, I^n)$, $\theta \in \mathbb{R}$ et que l'on veuille tester $H_0: \theta = 0$ contre $H_1: \theta = 1$ à partir d'un échantillon (X_1, \dots, X_n) du modèle $X_i \sim N(\theta, 1)$. Nous verrons qu'une telle loi X_n est une statistique de test pertinente. Par exemple, si $X_1(\omega) = 0.8$, que va-t-on choisir entre H_0 et H_1 ? Naturellement, une règle critique sera de la forme ϕ_s , $s \in \mathbb{R}$, car X_n est un estimateur de θ . On déterminera s à l'aide de $P_0(\text{Choisir } H_1) = \alpha = P_0(X_1 \geq s)$, donc par exemple, $\alpha = 0.05$, $s \approx 1.6$. Par suite si $X_1(\omega) = 0.8$, on accepte H_0 et l'erreur de seconde espèce $P_1(X_1 < s) \approx 0.74$, donc le test est : le test n'est pas discriminant. Maintenant, si on inverse H_0 et H_1 : $H_0: \theta = 1$ contre $H_1: \theta = 0$, le même ϕ_s avec $s = 1.6$, conduit à accepter H_1 avec une erreur de seconde espèce ≈ 0.74 . On obtient donc des résultats opposés pour la même expérience à priori. Les hypothèses H_0 et H_1 ne sont clairement pas interchangeables.

La question principale maintenant est de savoir comment trouver une statistique de test pertinente dans ce cadre pour décider d'utiliser un estimateur du maximum de vraisemblance.

4.2 Test de Wald

Un estimateur du maximum de vraisemblance permet d'associer à chaque observation x la même "forme" qu'un $\theta \in \Theta$. Cependant, il est difficile de trouver la bonne l'estimateur du maximum de vraisemblance θ . Si θ est fixé, il est possible, on utilisera directement θ comme statistique de test.

Sinon, de manière plus générale, on connaît la loi asymptotique de $\hat{\theta}_n$ quand le modèle est régulier. Donc quand n est grand, on pourra utiliser une loi normale comme approximation. Mais, nous devons vérifier la plausibilité de la même trajectoire de covariances asymptotique, qui est la matrice d'information de Fisher inverse dépend du paramètre θ . Aussi va-t-on préférer utiliser la statistique de test T suivante :

Définissons pour un modèle paramétrique donné $(\mathcal{Q}, \mathcal{A}_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$. La statistique de Wald T pour le test $H_0: \theta = \theta_0$ contre $H_1: \theta \in \Theta_1$ est $T_n = n \cdot l'(\theta_n - \theta_0) \cdot l(\theta_n - \theta_0)$.

Pour montrer "formellement" la pertinence de ce test, on considère la suite de tests T_n se plaçant dans le "grand" modèle asymptotique :

Théorème Dans le cadre d'un modèle paramétrique $((\mathcal{Q}^N, \mathcal{A}_{IN}, (f_\theta \cdot d\mu)^{\otimes N}, \theta \in \Theta)$, où $\theta \in \mathbb{R}^p$, donner par un échantillon θ_n pour le problème de test $H_0: \theta = \theta_0$ contre $H_1: \theta \in \Theta_1$, alors, en notant T_n la statistique de test de Wald pour le problème de taille n sous l'hypothèse H_0 ,

$$T_n \xrightarrow[n \rightarrow +\infty]{L} \chi^2(p).$$

La région de rejet asymptotique du test sera donc de la forme $T_n \geq t_n$ est le quantile d'ordre α de la loi $\chi^2(p)$. La suite de tests T_n a donc une puissance qui tend vers 1 lorsque α est fixé.

Démonstration : la loi asymptotique, déduit la loi asymptotique de $\sqrt{n} \cdot l(\theta_n - \theta_0)$ suit asymptotiquement une loi $N(0, I)$ sous l'hypothèse H_0 et $T_n = \sqrt{n} \cdot l(\theta_n - \theta_0)^2$. ■

Voici donc un premier type de test, sous certaines conditions de régularité du modèle et pour certaines hypothèses de tests existantes. Mais pourrions-nous faire mieux ? Et en réponse, oui ! Nous faut donc définir un moyen de comparer les deux tests.

4.3 Test du rapport de vraisemblance

Définissons les hypothèses et notations précédentes. On dira qu'un test est uniformément plus puissant (U.P.P.) au seuil α si le niveau de ϕ est inférieur à la statistique de test inférieure à α et si pour tout autre test ψ au même niveau α , $\forall \theta \in \Theta_1$,

$$IE_\theta(\phi) = 1 - IP_\theta(T \in W) \leq 1 - IP_\theta(T' \in W) = IE_\theta(\psi).$$

Définissons sous les hypothèses précédentes, si L est la vraisemblance, on appellera test du rapport de vraisemblance (test de Neyman - Pearson dans le cas d'hypothèses simples) un test de statistique T tel que

$$T = \frac{\sup_{\theta \in \Theta_0} L_{\theta}(X_1, \dots, X_n)}{\sup_{\theta \in \Theta} L_{\theta}(X_1, \dots, X_n)}.$$

La région critique W associée à un test est de la forme $W =]K[$ (donc si $T < K$, on rejette H_0).

Une des vertus du test du rapport de vraisemblance par rapport au test de Wald est qu'il est utile dans un modèle non régulier (mais la question de savoir, ou de la loi d'une fonctionnelle de ce test, demeure). De plus, la propriété suivante confirme l'efficacité de cette statistique de test:

Propriété (Principe de Lehman). Dans le cas du test de deux hypothèses simples, ou d'un test en fait $(\Theta \subset \mathbb{R})$, ce test est U.P.P. Dans le cas d'un test simple, il n'existe pas forcément de test U.P.P.

Démonstration ■

Enfin, un tel test pour un modèle régulier, va pouvoir être aidé par le théorème de la normalité asymptotique de l'estimateur du maximum de vraisemblance:

Théorème. Dans le cadre d'un modèle paramétrique $((\mathcal{Q}^N, \mathcal{A}_N, (f_{\theta} \cdot d\mu)^{\otimes N}, \theta \in \Theta)$, où $\mathcal{Q} = \mathbb{R}^p$, donné par une mesure régulière, pour le problème de test $H: \theta = \theta_0$ contre $H: \theta \neq \theta_0$, alors, en notant T_n la statistique du rapport de vraisemblance pour le test de taille n ,

$$-2 \log(T_n) \xrightarrow[n \rightarrow +\infty]{L} \chi^2(p).$$

La région de rejet asymptotique du test sera donc de la forme $W =]c_{\alpha}[$, où c_{α} est le quantile d'ordre $1 - \alpha$ de la loi du $\chi^2(p)$. La suite de tests (T_n) a donc une puissance qui tend vers 1 lorsque α est fixé.

Démonstration: la démonstration reprend un peu celle de la normalité asymptotique du maximum de vraisemblance. ■

Université Paris I, Paris - Sorbonne

Première Année Master M.A.E.F. 2005-2006

Statistique

Cours de statistique, novembre 2005

Examen de 2 h 00. Tout document ou calculatrice est interdit.

1. On considère une suite $(X_n)_{n \in \mathbb{N}}$ de variables aléatoires indépendantes et identiquement distribuées suivant une loi $\mathcal{N}(m, \sigma^2)$, où $m \in \mathbb{R}$ et $\sigma^2 > 0$ sont des paramètres connus.

Soit également $p \in \mathbb{N}^*$, $\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n)$ et $\sigma_n^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2$.

- Pour n fixé, déterminer le mode statistique par échantillon.
- Pour $n \in \mathbb{N}^*$, quelle est la loi de \bar{X}_n ?
- Quelles sont les limites (en probabilité) de \bar{X}_n et de σ_n^2 (justifier...)?
- Montrer que la convergence de (\bar{X}_n) induit celle de (X_1, \dots, X_n) . Déterminer la loi du vecteur (X_1, \dots, X_n) . Les (X_k) sont-elles indépendantes?
- Soit $\bar{\bar{X}}_n = \frac{1}{n}(\bar{X}_1 + \dots + \bar{X}_n)$. Quelle est la loi de $\bar{\bar{X}}_n$ pour $n \in \mathbb{N}^*$? En déduire que $\bar{\bar{X}}_n \xrightarrow[n \rightarrow +\infty]{p.s.} m$. Montrer également que $\bar{\bar{X}}_n \xrightarrow[n \rightarrow +\infty]{p.s.} m$.
- Comment peut-on faire pour savoir quelle suite de variables aléatoires $(\bar{X}_k)_{k \geq 1}$ s'approche le plus vite de m ? Conclusion?
- Pour le mode statistique de taille n^2 est-il supposable, montrer que la statistique \bar{X}_n est exhaustive pour $n \in \mathbb{N}^*$ et la statistique (X_1, \dots, X_n) ? Enfin, la statistique \bar{X}_n est-elle exhaustive?

2. En fait, on ne connaît pas explicitement que X_1 est plutôt pour tout $n \in \mathbb{N}^*$, $T_k = \max(X_1, \dots, X_k)$.

- La convergence de (T_n) induit-elle celle de (X_1, X_n) ?
- Déterminer la fonction de répartition de T_k , puis, après avoir montré l'existence, sa densité par rapport à la mesure de Lebesgue, le tout en fonction de la fonction de répartition F et de la densité f .
- Déterminer, en justifiant, le comportement asymptotique de (T_n) .
- Pour $k \in \mathbb{N}^*$, montrer que T_k et T_{k+1} ne sont pas indépendantes. Montrer que $\mathbb{P}(T_{k+1} = T_k) = \frac{k}{k+1}$. En déduire la mesure de probabilité variable $T - T_k$. La loi de probabilité de la variable T est-elle "continue"? Pourquoi?
- La statistique T est-elle exhaustive pour le mode statistique de taille n^2 où n est supposable? Et la statistique (T, T_n) ?

Université Paris I, Paris - Sorbonne

Première Année Master M.A.E.F. 2005-2006

Statistique

Cours de statistique 2006

Exam en de 2 h 00. Tout document ou calculatrice est interdit.

1. Soit la variable X qui suit une loi dont la densité par rapport à la mesure de Lebesgue sur $]0, 1]$ est, avec $\beta \in \mathbb{R}$:

$$f_X(x) = K \cdot x^\beta \quad \text{pour tout } x \in]0, 1],$$

- (a) Déterminer K en fonction de β , préciser quelle condition doit être vérifiée pour que f_X soit une densité, puis en déduire $IE(X)$ et $var(X)$, en précisant également des conditions sur β .
- (b) On suppose que la suite (X_n) est constituée de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X . Soit un échantillon (X_1, \dots, X_n) . On désire estimer β à partir de cet échantillon. Quel est le modèle statistique ? Montrer que ce modèle appartient à la famille exponentielle.
- (c) En déduire qu'il n'existe pas d'estimateurs sans biais efficaces de β .
- (d) Montrer que $q_n(X) = -\frac{1}{n} \sum_{i=1}^n \log(X_i)$ est un estimateur sans biais de β (utiliser les lois gamma...), puis qu'il est de variance minimale parmi les estimateurs sans biais (Lehman-Scheffé).

2. Soit Y une variable suivant une loi de Bernoulli de paramètre p et indépendante de X . On définit une variable Z de la manière suivante : si $Y = 1$, alors $Z = X$, et si $Y = 0$ alors $Z = 0$.

- (a) Montrer que Z suit une loi absolument continue par rapport à la mesure de Lebesgue sur $[-1, 1]$ et que sa densité est :

$$f_Z(z) = (\beta + 1) |x|^\beta \cdot \mathbb{1}_{x \in]0, 1]} + (1 - p) \cdot \mathbb{1}_{x \in [-1, 0]} \quad \text{pour tout } x \in [-1, 1].$$

Calculer $IE(Z)$ et $var(Z)$ (en précisant les conditions sur β).

- (b) On suppose que la suite (Z_n) est constituée de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que Z . Soit un échantillon (Z_1, \dots, Z_n) . On désire estimer (β, p) à partir de cet échantillon. Quel est le modèle statistique ? Montrer que ce modèle appartient à la famille exponentielle.
- (c) En déduire une statistique exhaustive minimale et complète. Déterminer la matrice d'information de Fisher de (β, p) et la borne de Cramér-Rao. Déterminer une fonction $g(\beta, p)$ que l'on ne peut estimer sans biais et de manière efficace.

- (d) Déterminer, après avoir montré l'unicité de l'estimateur $(\hat{\beta}_n, \hat{p}_n)$ du maximum de vraisemblance de (β, p) , si les estimateurs $\hat{\beta}_n$ et \hat{p}_n sont-ils indépendants ? Déterminer un théorème de la limite centrale relative à $(\hat{\beta}_n, \hat{p}_n)$. Est-ce un estimateur asymptotiquement efficace ?
- (e) Déterminer une région de confiance de niveau 95 % sur (β, p) , en utilisant 1/ l'estimateur efficace de $g(\beta, p)$ 2/ l'estimateur de maximum de vraisemblance dans le cadre asymptotique.

Première Année Master M.A.E.F. 2005-2006

Statistique

Examen terminal, janvier 2006

Exam en de 3 h 00 . Tout document ou calculatrice est interdit.

1. On considère $(X_k)_{k \in \mathbb{N}}$ et $(X'_k)_{k \in \mathbb{N}}$ deux suites indépendantes de variables aléatoires définies sur le même espace de probabilité indépendantes et identiquement distribuées suivant les lois respectives $N(\mu, \sigma^2)$ (pour les X) et $N(\mu', \sigma'^2)$ (pour les X'), où $(\mu, \mu') \in \mathbb{R}^2$ et $\sigma^2 > 0$. Le but du problème est de tester $H_0: \mu = \mu'$ à partir d'une échantillon de chacune de ces suites.

Soit (X_1, \dots, X_n) et $(X'_1, \dots, X'_{n'})$, où $n \in \mathbb{N}^*$ et $n' \in \mathbb{N}^*$, deux échantillons issus de $(X_k)_{k \in \mathbb{N}}$ et $(X'_k)_{k \in \mathbb{N}}$. On pose $Z = (Z_1, \dots, Z_{n+n'}) = (X_1, \dots, X_n, X'_1, \dots, X'_{n'})$

- (a) Déterminer le modèle statistique associé à $Z = (Z_1, \dots, Z_{n+n'})$.
 (b) Montrer que ce modèle est exponentiel. En déduire que μ peut être estimé efficacement (on notera les estimateurs respectifs).
 (c) Montrer qu'un estimateur du maximum de vraisemblance existe et que :

$$\hat{\sigma}^2 = \frac{1}{n+n'} \sum_{j=1}^n Z_j^2 - \hat{\mu}^2 + \frac{1}{n+n'} \sum_{j=n+1}^{n+n'} Z_j^2 - \hat{\mu}'^2.$$

Est-ce un estimateur efficace ? Est-il convergent ? Efficace ?

- (d) Lorsque n et n' sont "grandes", en déduire de ce qui précède, des intervalles de confiance à 95 % pour μ et μ' .
 (e) Soit le problème de test :

$$H_0 : \mu = \mu' \quad \text{contre} \quad H_1 : \mu \neq \mu'$$

(σ^2 restant inconnu). Montrer que la statistique T du rapport de vraisemblance vérifie :

$$T = \frac{\hat{\sigma}^2}{\bar{\sigma}^2}^{(n+n')/2} \quad \text{avec} \quad \begin{cases} \hat{\sigma}^2 = \frac{1}{n+n'} \sum_{j=1}^{n+n'} Z_j^2 - \bar{Z}_{n+n'}^2 \\ \bar{Z}_{n+n'} = \frac{1}{n+n'} \sum_{i=1}^{n+n'} Z_i \end{cases}$$

En déduire que la région d'acceptation du test est sous la forme $K \leq T \leq K'$, avec K dépendant du niveau du test.

(f) Pour déterminer la valeur K en fonction du niveau α , on peut considérer la statistique, dite de Fisher,

$$T' = \frac{(n+n') \cdot \bar{\sigma}^2 - (n+n') \cdot \sigma^2}{\frac{n+n'}{n+n'-2} \sigma^2}$$

- i. Soit les vecteurs de $\mathbb{R}^{n+n'}$ $u_1 = (1, \dots, 1)$, $u = (1, \dots, 1, 0, \dots, 0)$ (soit n fois 1 et n' fois 0) et $u = u_1 - u$. Montrer que u est orthogonal aux u_i .
- ii. Sous l'hypothèse H_0 , déterminer une expression plus simple de T' , projeté orthogonal de Z sur le sous-espace vectoriel $\langle u \rangle$, et $P_{\langle u, u \rangle}(Z)$, projeté orthogonal de Z sur le s.e.v. engendré par u .
- iii. Montrer que sous l'hypothèse H_0 , le vecteur Z peut s'écrire $Z = \mu u_1 + \sigma \cdot \varepsilon$, où ε est un vecteur aléatoire gaussien composé de $n+n'$ variables gaussiennes centrées réduites.
- iv. Montrer que sous l'hypothèse H_0 , $(n+n') \cdot \sigma^2 = \sigma^2 \cdot P_A(\varepsilon)^2$, où u_1 est la norme euclidienne classique sur $\mathbb{R}^{n+n'}$ et A est un s.e.v. de $\mathbb{R}^{n+n'}$ que vous préciserez.
- v. En utilisant le Théorème de Pythagore, montrer que sous l'hypothèse H_0 , $(n+n') \cdot \bar{\sigma}^2 - (n+n') \cdot \sigma^2 = \sigma^2 \cdot P_B(\varepsilon)^2$, où B est un s.e.v. de $\mathbb{R}^{n+n'}$ que vous préciserez.
- vi. En utilisant le Théorème de Cochran, montrer que sous l'hypothèse H_0 , T' suit une loi de Fisher à $(1, (n+n'-2))$ degrés de liberté et Lo rsque n et n' sont "grands", quelle loi suit approximativement T' ?
- vii. Pour finir, déterminer K en fonction d'un quantile de la loi de Fisher $F_{(1, (n+n'-2))}$ de degré de liberté $n+n'-2$.

2. Soit X une variable aléatoire dont la mesure de probabilité est absolument continue par rapport à la mesure de Lebesgue sur \mathbb{R} et de densité

$$f(x) = \frac{\lambda}{2} \exp(-\lambda|x-m|) \quad \text{pour } x \in \mathbb{R},$$

avec $m \in \mathbb{R}$ et $\lambda > 0$, des paramètres inconnus.

- (a) Calculer l'espérance et la variance de X .
- (b) Calculer $IP(X = m)$ et $IP(X < m)$. En déduire la médiane (théorique) de la loi de X .
- (c) Soit une suite $(X_i)_{i \in \mathbb{N}}$ de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X , dont on extrait un échantillon (X_1, \dots, X_{2n+1}) . Par ailleurs, on note $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(2n+1)}$ la statistique d'ordre associée. Soit :

$$H_n(a) = \frac{1}{2n+1} \sum_{i=1}^{2n+1} |X_i - a| \quad \text{pour } a \in \mathbb{R}.$$

Calculer $H(X_{(n+1)})$ en fonction des X_i . Montrer que la fonction $H_n(a)$ est minimale en $X_{(n+1)}$ (on pourra éventuellement utiliser que $H_n(X_{(n+k)})$ en fonction des X_i pour $k > 1$).

- (d) On suppose que $m = 1$, donc que m est connu ($\lambda > 0$ restant inconnu). Quel est alors le modèle statistique ? Montrer que ce modèle appartient à la famille exponentielle, et déduire une statistique exhaustive dont vous montrerez qu'elle est complète. Déterminer la matrice d'information de Fisher. Quelle est la fonction de λ (à une transformation affine près) qui peut estimer efficacement ? Déterminer l'estimateur de maximum de vraisemblance de λ et montrer qu'il vérifie un théorème de la limite centrale.
- (e) On suppose désormais que $\lambda \in \mathbb{R}$ est inconnu, tout comme $m > 0$. Quel est alors le modèle statistique ? Montrer que ce modèle appartient à la famille exponentielle. À l'aide de la question 2, déterminer un estimateur λ_n du maximum de vraisemblance du couple (m, λ) .
- (f) Pour $a \in \mathbb{R}$, démontrer que λ_n converge presque sûrement vers a vers $\mathbb{E}(|X - a|)$. Montrer que la fonction $a \mapsto \mathbb{E}(|X - a|)$ est minimale en $a = m$. En déduire que $m \xrightarrow[n \rightarrow +\infty]{p.s.} m$, puisque $\lambda_n \xrightarrow[n \rightarrow +\infty]{p.s.} \lambda$.

Première Année Master M.A.E.F. 2005-2006

Statistique

Examen de septembre 2006

Exam en de 3 h 00. Tout document ou calculatrice est interdit.

1. On considère $(X_k)_{k \in \mathbb{N}}$ une suite de variables aléatoires indépendantes et identiquement distribuées suivant la loi $\mathcal{N}(m, 1)$ où $m \in \mathbb{R}$. Soit la suite de variables $(Y_k)_{k \in \mathbb{N}}$ telle que pour tout $k \in \mathbb{N}^*$,

$$Y_k = X_1 + \dots + X_k.$$

Soit (Y_1, \dots, Y_n) , où $n \in \mathbb{N}^*$ un échantillon issu de $(Y_k)_{k \in \mathbb{N}}$.

- (a) Déterminer la loi de Y_k pour $k \in \mathbb{N}^*$.
 (b) Déterminer la loi du vecteur (Y_1, \dots, Y_n) . Montrer que pour tout j , Y_j n'est pas indépendante de Y .
 (c) Déterminer le modèle statistique associé à (Y_1, \dots, Y_n) .
 (d) Montrer que ce modèle est exponentiel.
 (e) Soit J_n la matrice de covariance de (Y_1, \dots, Y_n) . Vérifier que

$$J_n^{-1} = \begin{pmatrix} 2 & -1 & 0 & 0 & \dots & 0 & 0 \\ -1 & 2 & -1 & 0 & \dots & 0 & 0 \\ 0 & -1 & 2 & -1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & -1 & 1 \end{pmatrix}.$$

En déduire que m peut être estimé par un estimateur m_n (quel est-il) sans biais et efficace.

- (f) Déterminer l'estimateur du maximum de vraisemblance de m . Est-ce un estimateur sans biais ? Est-il convergent ? Efficace ? Quel est son risque quadratique ?
 (g) Déterminer la statistique du test de rapport de vraisemblance pour le test

$$H_0 : m = m_0 \quad \text{contre} \quad H_1 : m \neq m_0,$$

où m_0 est une constante connue. On utilise une application T_n de ce test au niveau 5 % pour $n = 100$. On trouve que $T_n(m_0) = 0.05$. Accepte-t-on alors l'hypothèse H_0 ?

2. Soit X une variable aléatoire dont la mesure de probabilité est absolument continue par rapport à la mesure de Lebesgue sur \mathbb{R} et de densité

$$f(x) = k \cdot \frac{1}{\sqrt{x}} \cdot \mathbb{I}_{0 \leq x \leq \theta} \quad \text{pour } x \in \mathbb{R},$$

avec $k \in \mathbb{R}$ et $\theta > 0$, des paramètres inconnus.

- (a) Déterminer l'expression de k en fonction de θ pour l'estimation de θ par la méthode des moments et la variance de X .
- (b) Soit une suite $(X_i)_{i \in \mathbb{N}}$ de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X , dont on extrait un échantillon (X_1, \dots, X_n) . Par ailleurs, on note $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ la statistique d'ordre associée. On désire estimer θ à partir de $(X_{(1)}, \dots, X_{(n)})$. Quel est alors le meilleur estimateur ? Quel est la vraisemblance pour ce modèle ?
- (c) Montrer que $T = X_{(n)} = \max(X_1, \dots, X_n)$ est une statistique exhaustive pour ce modèle.
- (d) Montrer que cette statistique est minimale.
- (e) Soit $(x_1, \dots, x_n) \in \mathbb{R}^n$ et le même n -uplet ordonné $(x) = (x_{(1)} \leq \dots \leq x_{(n)}) = \max(x)$. Montrer que :

$$\mathbb{P}(X_{(1)} \leq x_{(1)}, \dots, X_{(n)} \leq x_{(n)}) = n! \mathbb{P}(X_1 \leq x_{(1)} \cap X_1 \leq X_2 \leq x_{(2)} \cap \dots \cap X_{n-1} \leq X_n \leq x_{(n)}).$$

Montrer par récurrence sur n les ordres partielles que :

$$\frac{\partial^n}{\partial x_{(1)} \dots \partial x_{(n)}} \mathbb{P}(X_1 \leq x_{(1)} \cap X_1 \leq X_2 \leq x_{(2)} \cap \dots \cap X_{n-1} \leq X_n \leq x_{(n)}) = \prod_{i=1}^n f(x_{(i)}).$$

Soit $L_{\theta}^{(n)}$ la vraisemblance de $(X_{(1)}, \dots, X_{(n)})$. Démontrer de ce qui précède que :

$$L_{\theta}^{(n)}(x_{(1)}, \dots, x_{(n)}) = n! \cdot L_{\theta}(x_1, \dots, x_n).$$

- (f) Déterminer la densité puis le biais de T .
- (g) Montrer que T est une statistique exhaustive et complète.
- (h) Démontrer de ce qui précède un estimateur de θ , sans biais et uniformément de variance minimale.
- (i) Calculer le risque quadratique de T . Démontrer que $T \xrightarrow[n \rightarrow +\infty]{P} \theta$, puis, que pour tout $\alpha \in [0, 1[$, $\mathbb{P}(T_n - \theta) \xrightarrow[n \rightarrow +\infty]{P} 0$.
- (j) Déterminer explicitement un intervalle de confiance à 95 % de θ en fonction de n .

Université Paris I, Paris Lodron - Salzburg

Première Année Master M.A.E.F. 2006-2007

Statistique

Cours de statistique, novembre 2006

Examen de 2 h 00. Tout document ou calculatrice est interdit.

1. On considère une suite (ε_n) de variables aléatoires indépendantes et identiquement distribuées suivant une loi $\mathcal{N}(0, \sigma^2)$, où $\sigma^2 > 0$ est un paramètre inconnu. Pour tout $n \in \mathbb{N}^*$, on définit :

$$X_n = \varepsilon_n - \alpha \varepsilon_{n-1},$$

où $\alpha \in \mathbb{R}$ est inconnu. On notera par la suite σ_α^2 (σ

(a) Montrer que $\sigma_\alpha^2 \in 2\sigma^2$.

(b) Pour tout $n \in \mathbb{N}^*$, déterminer $E(X_n)$ et $\text{Var}(X_n)$. Montrer que (X_n) est une suite de variables indépendamment distribuées suivant la loi.

(c) Montrer que $\text{cov}(X_i, X_j) = (1 + \alpha^2)\sigma^2$ si $i = j$, $\text{cov}(X_i, X_j) = -\alpha\sigma^2$ si $|i - j| = 1$ et $\text{cov}(X_i, X_j) = 0$ sinon.

(d) Pour n fixé, en déduire la loi du vecteur (X_1, \dots, X_n) , puis déterminer le module statistique par la méthode appropriée en précisant une mesure de convergence.

(e) Soit $\hat{\alpha}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2$. Montrer que $\hat{\alpha}_n^2$ est un estimateur non biaisé. Soit

$$Z_n^{(1)} = \frac{1}{n} \sum_{k=1}^{[n/2]} X_{2k}^2 \quad \text{et} \quad Z_n^{(2)} = \frac{1}{n} \sum_{k=1}^{[(n+1)/2]} X_{2k-1}^2$$

avec $[x]$ la partie entière de x . Montrer que les suites de variables aléatoires $(Z_n^{(1)})_n$ convergent presque sûrement (ou presque partout) et qu'elles vérifient un théorème de la limite centrale.

(f) Montrer que si deux suites de variables aléatoires convergent presque sûrement, la suite posée de leurs sommes converge presque sûrement. En déduire que $\hat{\alpha}_n^2$ converge presque sûrement vers σ_α^2 .

(g) Soit $\hat{\alpha}_n = \frac{1}{n-1} \sum_{i=1}^{n-1} X_i X_{i+1}$. En utilisant le même type d'argument que précédemment, montrer que $\hat{\alpha}_n$ converge presque sûrement vers α . En déduire un estimateur de α convergent presque sûrement.

2. Soit X une variable aléatoire dont la loi est absolument continue par rapport à la mesure de Lebesgue sur \mathbb{R} et telle que sa densité appartienne à cette mesure soit :

$$f_{\alpha,a}(x) = K \cdot \frac{1}{|x|^\alpha} \mathbb{I}_{-a \leq x \leq a}$$

avec $a \in \mathbb{R}$ et $\alpha \in \mathbb{R}$. Soit (X_1, \dots, X_n) un échantillon de n v.a. i.i.d. de loi θ que

- Après avoir précisé l'ensemble des valeurs Θ pour $\theta = (\alpha, \beta)$, déterminer K en fonction de a et α .
- Calculer $E(X)$ et vérifier que ce calcul peut être effectué pour $\theta \in \Theta$.
- Quel est le modèle statistique associé à (X_1, \dots, X_n) ?
- Montrer que $S = (X_1, \dots, X_n)$ est une statistique exhaustive. Montrer que pour $n \geq 3$ cette statistique n'est pas minimale.
- Déterminer une statistique $T, T_2(X)$ à valeurs dans \mathbb{R}^2 qui soit exhaustive minimale pour tout $\theta \in \Theta$.

Université Paris I, Paris - Sorbonne

Première Année Master M.A.E.F. 2006-2007

Statistique

Cours de statistique, janvier 2007

Exam en de 2 h 00. Tout document ou calculatrice est interdit.

1. Soit X une variable aléatoire suivant la loi suivante:

$$P(X = 1) = P(Y = -1) = p \text{ et } P(X = 0) = 1 - 2p,$$

où p est un paramètre inconnu.

(a) Déterminer l'ensemble des valeurs possibles pour p . Calculer EX et $\text{var}X$.

(b) On suppose que la suite $(X_i)_{i \in \mathbb{N}}$ est constituée de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X . Soit un échantillon (X_1, \dots, X_n) . Déterminer le maximum de la statistique associée à cet échantillon et déterminer une mesure dominante de la famille. Montrer que le maximum n'appartient pas à la famille exponentielle. Déduire une statistique exhaustive et proposer un estimateur. Montrer que p n'est pas efficace et donner un tel estimateur. Calculer la borne de Cramér-Rao et vérifier qu'elle est bien atteinte par cet estimateur.

(c) On définit la suite $(Y_i)_{i \in \mathbb{N}^*}$ à partir de $(X_i)_{i \in \mathbb{N}}$ de la manière suivante:

$$Y_{i+1} = X_i \cdot X_{i+1} \text{ pour } i \in \mathbb{N}.$$

Déterminer la loi de Y_1 . Montrer que $\text{cov}(Y_i, Y_{i+1}) = 0$. Les (Y_i) sont-elles indépendantes?

(d) Montrer que $(|Y_1|, \dots, |Y_n|)$ est une statistique exhaustive pour la loi. Montrer que la statistique induit par (Y_1, \dots, Y_n) .

2. Soit la variable X qui suit une loi dont la densité rapportée à la mesure de Lebesgue sur \mathbb{R} est, avec $\theta > 0$ et $\alpha > 0$:

$$f_X(x) = K \cdot x^\alpha \mathbb{1}_{0 \leq x \leq \theta} \text{ pour tout } x \in \mathbb{R},$$

(a) Déterminer K en fonction de α et θ .

(b) Montrer que $Y = \log(\theta/X)$ suit une loi exponentielle d'ordre 1.

(c) On suppose que la suite $(X_i)_{i \in \mathbb{N}}$ est constituée de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X . Soit un échantillon (X_1, \dots, X_n) . On suppose que (θ, α) est inconnu. Déterminer alors le maximum de la statistique associée à cet échantillon et la mesure dominante. Calculer la borne de Cramér-Rao et vérifier qu'elle est bien atteinte par cet estimateur?

- (d) Dans cette question, et uniquement dans cette question, on suppose que θ est connu. Préciser alors le modèle statistique que le modèle appartient à la famille exponentielle ? Montrer que l'estimateur du maximum de vraisemblance existe, est unique et écrire :

$$\hat{\alpha} = \frac{1}{\frac{1}{n} \sum_{i=1}^n \log(\theta / X_i)} - 1$$

Montrer que $\hat{\alpha}$ converge presque sûrement vers α et qu'un théorème de la limite centrale que l'on peut en déduire un intervalle de confiance 95% sur α pour n grand.

- (e) Dans cette question, θ et α sont inconnus. Déterminer une statistique exhaustive pour le modèle. En vous aidant de la question précédente, déterminer l'estimateur du maximum de vraisemblance de (θ, α) et déterminer la fonction de répartition de $\log(\theta / X)$. En déduire que $\theta \xrightarrow[n \rightarrow +\infty]{P} \theta$, puis que $\sqrt{n} \log(\theta / X) \xrightarrow[n \rightarrow +\infty]{P} 0$.

- (f) Soit $(U_n)_{n \in \mathbb{N}}$ et $(V_n)_{n \in \mathbb{N}}$ deux suites de variables aléatoires définies sur le même espace de probabilité. Montrer que si (U_n) converge vers une loi P et (V_n) converge en probabilité vers 0, alors $(U_n + V_n)$ converge en loi vers P . On pourra par exemple majorer la différence de fonctions caractéristiques. En déduire que suit le même théorème de la limite centrale que $\hat{\alpha}$.

Université Paris I, Paris - Sorbonne

Première Année Master M.A.E.F. 2006-2007

Statistiques

Examen terminal, janvier 2007

Examen de 3 h 00. Tout document ou calculatrice est interdit.

1. On considère une suite de variables aléatoires $(X_k)_{k \in \mathbb{N}^*}$ définies sur le même espace de probabilité indépendantes et identiquement distribuées suivant une loi de Bernoulli de paramètre p . On définit :

$$Y = \min \{k \in \mathbb{N}^*, X_k = 0\}.$$

- (a) Comment peut-on interpréter la variable Y ? Montrer que la loi de Y est :

$$P(Y = k) = p^{k-1}(1-p) \quad \text{pour } k \in \mathbb{N}^*.$$

- (b) On suppose que (Y, X) est un échantillon de variables aléatoires indépendantes et identiquement distribuées suivant la loi Y avec $p \in]0, 1[$. Est-il intéressant de déterminer le mode statistique et sa mesure de minimum ? Montrer que ce mode est exponentiel. Ensuite, un estimateur sans biais et efficace de p est-il existant ? Déterminer la borne de Cramer-Rao et vérifier que cette borne est bien atteinte par

- (c) La variable Y est trop grande lorsqu'elle est trop "grande", par exemple, c'est-à-dire que l'on définit une variable telle que $Y = \min(Y, T)$.

2. Soit X une variable aléatoire dont la mesure de probabilité est absolument continue par rapport à la mesure de Lebesgue sur \mathbb{R} et de densité

$$f(x) = \frac{\lambda}{2} \exp(-\lambda|x - m|) \quad \text{pour } x \in \mathbb{R},$$

avec $m \in \mathbb{R}$ et $\lambda > 0$, des paramètres inconnus.

- (a) Calculer l'espérance et la variance de X .

- (b) Calculer $P(X = m)$ et $P(X < m)$. En déduire la médiane (théorique) de la loi de X .

- (c) Soit une suite $(X_i)_{i \in \mathbb{N}}$ de variables aléatoires indépendantes et identiquement distribuées suivant la même loi que X , dont on extrait un échantillon (X_1, \dots, X_{2n+1}) . Par ailleurs, on note $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(2n+1)}$ la statistique d'ordre associée. Soit :

$$H_n(a) = \frac{1}{2^{n+1}} \sum_{i=1}^{2n+1} |X_i - a| \quad \text{pour } a \in \mathbb{R}.$$

Calculer $H(X_{(n+1)})$ en fonction de n . Montrer que la fonction $H_n(a)$ est minimale en $X_{(n+1)}$ (on pourra éventuellement utiliser que $H(X_{(n+k)})$ en fonction de n pour $k > 1$).

- (d) On suppose que $m = 1$, donc que m est connu ($\lambda > 0$ restant inconnu). Quel est alors le mode statistique ? Montrer que ce mode appartient à la famille exponentielle, et déduire une statistique exhaustive dont vous montrerez qu'elle est complète. Déterminer la matrice d'information de Fisher. Quelle est la fonction de λ (à une transformation affine près) qui est le plus efficace ? Déterminer l'estimateur de maximum de vraisemblance de λ et montrer qu'il vérifie un théorème de la limite centrale.
- (e) On suppose désormais que $\lambda \in \mathbb{R}$ est inconnu, tout comme $m > 0$. Quel est alors le mode statistique ? Montrer que ce mode appartient à la famille exponentielle. À l'aide de la question 2, déterminer un estimateur $(\hat{m}, \hat{\lambda}_n)$ du maximum de vraisemblance du couple (m, λ) .
- (f) Pour $a \in \mathbb{R}$, démontrer que (\hat{a}) converge presque sûrement vers a vers $IE(|X - a|)$. Montrer que la fonction $a \mapsto IE(|X - a|)$ est minimale en $a = m$. En déduire que $\hat{m} \xrightarrow[n \rightarrow +\infty]{p.s.} m$, puisque $\lambda_n \xrightarrow[n \rightarrow +\infty]{p.s.} \lambda$.

Première Année Master M.A.E.F. 2006-2007

Statistique

Examen de septembre 2007

Examen de 3 h 00. Tout document ou calculatrice est interdit.

1. Soit la fonction $f(x) = \frac{1}{2}(a - a^2 \cdot x) \cdot \mathbb{I}_{\{-1/a \leq x \leq 1/a\}}$ où $a > 0$.

- (a) Montrer que f est une densité de probabilité par rapport à la mesure de Lebesgue et la tracer.
- (b) On suppose que X est une variable aléatoire de densité f . Déterminer EX et $\text{var}X$.
- (c) Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires indépendantes telle que la densité de X_n soit f_n pour tout $n \in \mathbb{N}^*$. Déterminer la limite en probabilité de $(X_n)_{n \in \mathbb{N}}$ lorsque $n \rightarrow \infty$.
- (d) Soit (X_1, \dots, X_n) un échantillon de v.a. i.i.d. de densité f . On suppose que a est inconnu. Déterminer l'estimateur du maximum de vraisemblance de a . Calculer la fonction de répartition de a et en déduire sa convergence en probabilité vers a .
- (e) Pour $\alpha > 0$, déterminer un intervalle de confiance de niveau α pour a .
- (f) Déterminer le test du rapport de vraisemblance de niveau α pour tester l'hypothèse $H_0 : a = a_0$, contre l'hypothèse $H_1 : a = a_1$ et déterminer la zone d'acceptation du test en fonction de α .
- (g) Proposer un autre estimateur convergent de a .

2. Une compagnie fabrique des piles et s'intéresse à savoir quelle est leur durée moyenne. Pour ce faire, on considère 10 000 piles produites chaque jour que l'on soumet à la mise à l'épreuve. Comme on ne veut pas attendre que toutes les piles soient usées, on décide d'arrêter l'expérience au bout de 10 jours et de compter combien sont encore "en vie". Soit N le nombre.

- (a) Dans une première approximation, on suppose que la durée d'une pile peut être modélisée par une loi exponentielle de paramètre $\lambda > 0$. Quelle est la loi de la durée de vie moyenne (arithmétique) T d'une pile en fonction de λ ? Quelle est, en fonction de T , la probabilité qu'une pile meure avant 10 jours? Montrer que $N/10 000$ suit approximativement une loi normale de la limite centrale dont on peut alors se servir en fonction de N pour déduire la loi d'un estimateur T_N de T en fonction de $N/10 000$ dont on donnera la loi de la limite centrale.
- (b) Montrer que N n'est pas une statistique exhaustive pour le paramètre λ par rapport à l'échantillon des 10 000 durées de vie des piles. Et si l'on avait attendu x jours au lieu de 10 ? Déterminer la loi de l'estimateur T_N par rapport à x et que T est une

mieux" T (do nc tro uver x q ue l so it de va ria nce min im a le) Ce ré sultat vo us
 semble t- il en pr a tiq ue é s s a nt?